

Real and Abstract Analysis

Kenneth Kuttler klkuttler@gmail.com

February 12, 2024

Contents

1	Review of Some Linear Algebra	15
1.1	The Matrix of a Linear Map	15
1.2	Block Multiplication of Matrices	15
1.3	Schur's Theorem	18
1.4	Hermitian and Symmetric Matrices	20
1.5	The Right Polar Factorization	21
1.6	Elementary matrices	24
1.7	The Row Reduced Echelon Form Of A Matrix	32
1.8	Finding the Inverse of a Matrix	35
1.9	The Mathematical Theory of Determinants	39
1.9.1	The Function sgn	39
1.9.2	The Definition of the Determinant	41
1.9.3	A Symmetric Definition	43
1.9.4	Basic Properties of the Determinant	44
1.9.5	Expansion Using Cofactors	46
1.9.6	A Formula for the Inverse	48
1.9.7	Cramer's Rule	49
1.9.8	Rank of a Matrix	50
1.9.9	An Identity of Cauchy	52
1.10	The Cayley Hamilton Theorem	53
I	Topology, Continuity, Algebra, Derivatives	55
2	Some Basic Topics	57
2.1	Basic Definitions	57
2.2	The Schroder Bernstein Theorem	59
2.3	Equivalence Relations	63
2.4	\sup and \inf	63
2.5	Double Series	64
2.6	$\lim \sup$ and $\lim \inf$	66
2.7	Nested Interval Lemma	68
2.8	The Hausdorff Maximal Theorem	68
3	Metric Spaces	71
3.1	Open and Closed Sets, Sequences, Limit Points	71
3.2	Cauchy Sequences, Completeness	73

3.3	Closure of a Set	74
3.4	Separable Metric Spaces	75
3.5	Compact Sets	76
3.6	Continuous Functions	80
3.7	Continuity and Compactness	81
3.8	Lipschitz Continuity and Contraction Maps	82
3.9	Convergence of Functions	84
3.10	Compactness in $C(X, Y)$ Ascoli Arzela Theorem	85
3.11	Connected Sets	88
3.12	Partitions of Unity in Metric Space	91
3.13	Completion of Metric Spaces	93
3.14	Exercises	94
4	Linear Spaces	99
4.1	Algebra in \mathbb{F}^n , Vector Spaces	99
4.2	Subspaces Spans and Bases	100
4.3	Inner Product and Normed Linear Spaces	103
4.3.1	The Inner Product in \mathbb{F}^n	103
4.3.2	General Inner Product Spaces	104
4.3.3	Normed Vector Spaces	106
4.3.4	The p Norms	106
4.3.5	Orthonormal Bases	108
4.4	Equivalence of Norms	109
4.5	Covering Theorems	113
4.5.1	Vitali Covering Theorem	113
4.5.2	Besicovitch Covering Theorem	115
4.6	Exercises	120
5	Functions on Normed Linear Spaces	125
5.1	$\mathcal{L}(V, W)$ as a Vector Space	125
5.2	The Norm of a Linear Map, Operator Norm	126
5.3	Comparisons	128
5.4	Continuous Functions in Normed Linear Space	129
5.5	Polynomials	130
5.6	Weierstrass Approximation Theorem	130
5.7	Functions of Many Variables	132
5.8	A Generalization	134
5.9	An Approach to the Integral	137
5.10	The Stone Weierstrass Approximation Theorem	141
5.11	Connectedness in Normed Linear Space	144
5.12	Saddle Points*	145
5.13	Exercises	150
6	Fixed Point Theorems	155
6.1	Simplices and Triangulations	155
6.2	Labeling Vertices	159
6.3	The Brouwer Fixed Point Theorem	161
6.4	The Schauder Fixed Point Theorem	163

6.5	The Kakutani Fixed Point Theorem	167
6.6	Ekeland's Variational Principle	170
6.6.1	Cariste Fixed Point Theorem	172
6.6.2	A Density Result	173
6.7	Exercises	175
7	The Derivative	183
7.1	Limits of a Function	183
7.2	Basic Definitions	186
7.3	The Chain Rule	188
7.4	The Matrix of the Derivative	188
7.5	A Mean Value Inequality	190
7.6	Existence of the Derivative, C^1 Functions	191
7.7	Higher Order Derivatives	193
7.8	Some Standard Notation	194
7.9	The Derivative and the Cartesian Product	195
7.10	Mixed Partial Derivatives	198
7.11	A Cofactor Identity	200
7.12	Newton's Method	201
7.13	Exercises	202
8	Implicit Function Theorem	205
8.1	Statement and Proof of the Theorem	205
8.2	More Derivatives	211
8.3	The Case of \mathbb{R}^n	212
8.4	Exercises	213
8.5	The Method of Lagrange Multipliers	216
8.6	The Taylor Formula	217
8.7	Second Derivative Test	218
8.8	The Rank Theorem	220
8.9	The Local Structure of C^1 Mappings	224
8.10	Invariance of Domain	227
8.11	Exercises	231
II	Integration	235
9	Measures and Measurable Functions	237
9.1	Simple Functions and Measurable Functions	237
9.2	Measures and their Properties	241
9.3	Dynkin's Lemma	243
9.4	Outer Measures	244
9.5	Measures From Outer Measures	245
9.6	Measurable Sets Include Borel Sets?	248
9.7	An Outer Measure on $\mathcal{P}(\mathbb{R})$	250
9.8	Measures and Regularity	252
9.9	One Dimensional Lebesgue Stieltjes Measure	257
9.10	Exercises	258

9.11	Completion of a Measure Space	260
9.12	Vitali Coverings	262
9.13	Differentiation of Increasing Functions	265
9.14	Exercises	268
9.15	Multifunctions and Their Measurability	271
9.15.1	The General Case	271
9.15.2	A Special Case When $\Gamma(\omega)$ Compact	273
9.15.3	Kuratowski's Theorem	273
9.15.4	Measurability of Fixed Points	275
9.15.5	Other Measurability Considerations	276
9.16	Exercises	278
10	The Abstract Lebesgue Integral	279
10.1	Nonnegative Measurable Functions	279
10.1.1	Riemann Integrals for Decreasing Functions	279
10.1.2	The Lebesgue Integral for Nonnegative Functions	280
10.2	Nonnegative Simple Functions	281
10.3	The Monotone Convergence Theorem	282
10.4	Other Definitions	283
10.5	Fatou's Lemma	283
10.6	The Integral's Righteous Algebraic Desires	284
10.7	The Lebesgue Integral, L^1	284
10.8	The Dominated Convergence Theorem	288
10.9	Some Important General Theory	291
10.9.1	Egoroff's Theorem	291
10.9.2	The Vitali Convergence Theorem	292
10.10	One Dimensional Lebesgue Stieltjes Integral	294
10.11	The Distribution Function	297
10.12	Good Lambda Inequality	299
10.13	Radon Nikodym Theorem	300
10.14	Iterated Integrals	304
10.15	Jensen's Inequality	309
10.16	Faddeyev's Lemma	310
10.17	Exercises	310
11	Regular Measures	315
11.1	Regular Measures in a Metric Space	315
11.2	Constructing Measures from Functionals	317
11.3	The p Dimensional Lebesgue Measure	320
11.4	Maximal Functions	323
11.5	Strong Estimates for Maximal Function	326
11.6	The Brouwer Fixed Point Theorem	327
11.7	Change of Variables, Linear Maps	329
11.8	Differentiable Functions and Measurability	331
11.9	Change of Variables, Nonlinear Maps	333
11.10	Mappings Not One to One	337
11.11	Spherical Coordinates	339
11.12	Symmetric Derivative for Radon Measures	341

11.13	Radon Nikodym Theorem, Radon Measures	343
11.14	Absolutely Continuous Functions	344
11.15	Total Variation	348
11.16	Exercises	350
12	The L^p Spaces	357
12.1	Basic Inequalities and Properties	357
12.2	Density Considerations	362
12.3	Separability	363
12.4	Continuity of Translation	365
12.5	Mollifiers and Density of Smooth Functions	366
12.6	Smooth Partitions of Unity	369
12.7	Exercises	371
13	Fourier Transforms	375
13.1	Fourier Transforms of Functions in \mathcal{G}	375
13.2	Fourier Transforms of Just about Anything	378
13.2.1	Fourier Transforms of Functions in $L^1(\mathbb{R}^n)$	381
13.2.2	Fourier Transforms of Functions in $L^2(\mathbb{R}^n)$	383
13.2.3	The Schwartz Class	386
13.2.4	Convolution	388
13.3	Exercises	390
14	Integration on Manifolds	393
14.1	Manifolds	393
14.2	The Area Measure on a Manifold	397
14.3	Divergence Theorem	400
14.4	Volumes of Balls in \mathbb{R}^p	404
14.5	Exercises	405
15	Degree Theory	409
15.1	Sard's Lemma and Approximation	411
15.2	Properties of the Degree	418
15.3	Brouwer Fixed Point Theorem	420
15.4	Borsuk's Theorem	421
15.5	Some Applications	424
15.6	Product Formula, Separation Theorem	426
15.7	General Jordan Separation Theorem	431
15.8	Uniqueness of the Degree	432
15.9	Exercises	433
16	Hausdorff Measure	437
16.1	Lipschitz Functions	437
16.2	Lipschitz Functions and Gateaux Derivatives	437
16.3	Rademacher's Theorem	438
16.4	Weak Derivatives	443
16.5	Definition of Hausdorff Measures	445
16.6	Properties of Hausdorff Measure	446
16.7	\mathcal{H}^p and m_p	447

16.8	Technical Considerations	449
16.8.1	Steiner Symmetrization	450
16.8.2	The Isodiametric Inequality	451
16.9	The Proper Value of $\beta(p)$	451
16.10	A Formula for $\alpha(p)$	452
17	The Area Formula	453
17.1	Estimates for Hausdorff Measure	453
17.2	Comparison Theorems	455
17.3	The Area Formula	456
17.4	The Divergence Theorem	461
17.5	The Coarea Formula	462
17.6	Change of Variables	467
17.7	Integration and the Degree	468
18	Differential Forms	471
18.1	The Wedge Product	475
18.2	The Exterior Derivative	476
18.3	Stokes Theorem	478
18.4	Lipschitz Maps	482
18.5	What Does it Mean?	483
18.6	Examples of $r([a, b])$	486
18.7	Orientation and Degree	486
18.8	Examples of Stoke's Theorem	488
18.8.1	Fundamental Theorem of Calculus	488
18.8.2	Line Integrals for Conservative Fields	488
18.8.3	Green's Theorem	488
18.8.4	Stoke's Theorem from Calculus	489
18.8.5	The Divergence Theorem	491
18.9	The Reynolds Transport Formula	492
18.10	Exercises	494
III	Abstract Theory	499
19	Hausdorff Spaces and Measures	501
19.1	General Topological Spaces	501
19.2	The Alexander Sub-basis Theorem	507
19.3	The Product Topology and Compactness	508
19.4	Stone Weierstrass Theorem	509
19.4.1	The Case of Locally Compact Sets	509
19.4.2	The Case of Complex Valued Functions	510
19.5	Partitions of Unity	511
19.6	Measures on Hausdorff Spaces	512
19.7	Measures and Positive Linear Functionals	516
19.8	Slicing Measures	520
19.9	Exercises	521
20	Product Measures	523

20.1	Algebras	523
20.2	Caratheodory Extension Theorem	524
20.3	Kolmogorov Extension Theorem	526
20.4	Exercises	530
21	Banach Spaces	533
21.1	Theorems Based on Baire Category	533
21.1.1	Baire Category Theorem	533
21.1.2	Uniform Boundedness Theorem	536
21.1.3	Open Mapping Theorem	536
21.1.4	Closed Graph Theorem	538
21.2	Hahn Banach Theorem	540
21.2.1	Partially Ordered Sets	540
21.2.2	Gauge Functions and Hahn Banach Theorem	540
21.2.3	The Complex Version of the Hahn Banach Theorem	542
21.2.4	The Dual Space and Adjoint Operators	543
21.3	Uniform Convexity of L^p	546
21.4	Closed Subspaces	552
21.5	Weak And Weak * Topologies	554
21.5.1	Basic Definitions	554
21.5.2	Banach Alaoglu Theorem	556
21.5.3	Eberlein Smulian Theorem	558
21.6	Differential Equations	561
21.7	Lyapunov Schmidt Procedure	563
21.8	The Holder Spaces	567
21.9	Exercises	569
22	Hilbert Spaces	575
22.1	Basic Theory	575
22.2	The Hilbert Space $L(U)$	580
22.3	Approximations in Hilbert Space	582
22.4	Orthonormal Sets	584
22.5	Compact Operators in Hilbert Space	586
22.5.1	Nuclear Operators	590
22.5.2	Hilbert Schmidt Operators	592
22.6	Roots of Positive Linear Maps	595
22.6.1	The Product of Positive Self Adjoint Operators	595
22.6.2	Roots of Positive Self Adjoint Operators	596
22.7	Differential Equations in Banach Space	599
22.8	General Theory of Continuous Semigroups	602
22.8.1	Generators of Semigroups	602
22.8.2	Hille Yosida Theorem	606
22.8.3	An Evolution Equation	610
22.8.4	Adjoint for Closed Operators, Hilbert Space	612
22.8.5	Adjoint, Reflexive Banach Space	615
22.9	Exercises	619
23	Representation Theorems	621

23.1	Radon Nikodym Theorem	621
23.2	Vector Measures	621
23.3	Representation for the Dual Space of L^p	626
23.4	Weak Compactness	631
23.5	The Dual Space of $L^\infty(\Omega)$	633
23.6	Non σ Finite Case	636
23.7	The Dual Space of $C_0(X)$	639
23.7.1	Extending Righteous Functionals	641
23.7.2	The Riesz Representation Theorem	642
23.8	Exercises	643
24	The Bochner Integral	647
24.1	Strong and Weak Measurability	647
24.1.1	Eggoroff's Theorem	654
24.2	The Bochner Integral	655
24.2.1	Definition and Basic Properties	655
24.2.2	Taking a Closed Operator Out of the Integral	659
24.3	Operator Valued Functions	662
24.3.1	Review of Hilbert Schmidt Theorem	663
24.3.2	Measurable Compact Operators	667
24.4	Fubini's Theorem for Bochner Integrals	668
24.5	The Spaces $L^p(\Omega; X)$	671
24.6	Measurable Representatives	677
24.7	Vector Measures	678
24.8	The Riesz Representation Theorem	682
24.9	An Example of Polish Space	687
24.10	Weakly Convergent Sequences	689
24.11	Some Embedding Theorems	690
24.12	Conditional Expectation in Banach Spaces	701
24.13	Exercises	704
25	Stone's Theorem and Partitions of Unity	705
25.1	Partitions of Unity and Stone's Theorem	708
25.2	An Extension Theorem, Retracts	710
IV	Stochastic Processes and Probability	713
26	Independence	715
26.1	Random Variables and Independence	715
26.2	Convergence in Probability	718
26.3	Kolmogorov Extension Theorem	719
26.4	Independent Events and σ Algebras	720
26.5	Banach Space Valued Random Variables	724
26.6	Reduction to Finite Dimensions	727
26.7	0, 1 Laws	727
26.8	Strong Law of Large Numbers	730
27	Analytical Considerations	735

27.1	The Characteristic Function	735
27.2	Conditional Probability	736
27.3	Conditional Expectation, Sub-martingales	740
27.4	Characteristic Functions and Independence	744
27.5	Characteristic Functions for Measures	748
27.6	Independence in Banach Space	750
27.7	Convolution and Sums	752
28	The Normal Distribution	759
28.1	The Multivariate Normal Distribution	759
28.2	Linear Combinations	762
28.3	Finding Moments	765
28.4	Prokhorov and Levy Theorems	766
28.5	The Central Limit Theorem	776
29	Martingales	781
29.1	Conditional Expectation	781
29.2	Conditional Expectation and Independence	784
29.3	Discrete Stochastic Processes	786
29.3.1	Upcrossings	788
29.3.2	The Sub-martingale Convergence Theorem	790
29.3.3	Doob Sub-martingale Estimates	793
29.4	Optional Sampling and Stopping Times	795
29.4.1	Optional Sampling for Martingales	798
29.4.2	Optional Sampling Theorem for Sub-Martingales	799
29.5	Reverse Sub-martingale Convergence Theorem	802
29.6	Strong Law of Large Numbers	804
30	Continuous Stochastic Processes	807
30.1	Fundamental Definitions and Properties	807
30.2	Kolmogorov Čentsov Continuity Theorem	809
30.3	Filtrations	817
30.4	Martingales and Sub-Martingales	825
30.5	Some Maximal Estimates	826
31	Optional Sampling Theorems	831
31.1	Review of Discreet Stopping Times	831
31.2	Review of Doob Optional Sampling Theorem	833
31.3	Doob Optional Sampling Continuous Case	834
31.3.1	Stopping Times	834
31.3.2	The Optional Sampling Theorem Continuous Case	839
31.4	Maximal Inequalities and Stopping Times	845
31.5	Continuous Sub-martingale Convergence	849
31.6	Hitting This Before That	853
31.7	The Space $\mathcal{M}_T^p(E)$	856
32	Quadratic Variation	859
32.1	How to Recognize a Martingale	859
32.2	Martingales and Total Variation	862

32.3	The Quadratic Variation	864
32.4	The Covariation	871
32.5	The Burkholder Davis Gundy Inequality	874
32.6	Approximation With Step Functions	880
33	Quadratic Variation and Stochastic Integration	883
33.1	The Stieltjes Integral	891
33.2	The Stochastic Integral When $f(s) \in \mathcal{L}_2(U, H)$	893
33.3	More on Stopping Times	898
33.4	Local Martingales as Integrators	901
33.5	The Stochastic Integral and the Quadratic Variation	903
33.6	The Case of $f \in L^2(\Omega \times [0, T]; \mathcal{L}_2(U, H))$	904

Copyright © 2018, You are welcome to use this, including copying it for use in classes or referring to it on line but not to publish it for money. I do constantly upgrade it when I find errors or things which could be improved, and I am constantly finding such things.

Preface

This book is on real and abstract analysis. There are four parts. The last part is an introduction to probability and stochastic processes. A course in multi-variable advanced calculus is contained in the first two. There is more there than one would have time to do and there are some advanced topics which are there because it is a convenient setting for them. Originally the first three parts were written for a multi-variable advanced calculus course which changed over time by inclusion of more advanced topics. However, it was easier to keep track of a single file. The emphasis is on finite dimensional spaces although not totally. There are things which are almost never discussed in multi-variable advanced calculus like the fixed point theorems. However, I think the Brouwer fixed point theorem is one of the most important theorems in mathematics and is being neglected along with its corollaries. I give several proofs in the exercises and in the book. There is too much reliance on these theorems without ever considering proofs. That is why there is a chapter on fixed point theorems. In general, I am trying to include all proofs of theorems instead of making the reader chase after them somewhere else or accept them on faith. I object to attempts to make mathematics functionally like a religion where we are asked to believe the words of authority figures. Of course, when you try to include all the proofs, you run the risk of making mistakes, and I certainly make my share, but one should at least try, even though it also results in a longer book.

I am reviewing a few topics from linear algebra mainly to refresh the memory or to read as needed, but I am assuming that people have had a reasonable course in linear algebra. Linear algebra should come before a book like this one.

I sometimes present important ideas more than once. Sometimes there is a special case in exercises and later the topic is discussed in the text. I think this can be useful in understanding some of these big theorems. Such duplication may not have been deliberate to begin with, but I have chosen to leave it in many cases.

Finite dimensional degree theory is neglected so there is a chapter on this also, presented as a part of analysis. It seems like it is common to neglect to give a careful treatment of the degree in \mathbb{R}^p . This is too bad. You end up missing out on fantastic finite dimensional topology like the Jordan separation theorem. I don't know a good proof for this without something like degree theory. Other somewhat unusual items are things like the Besicovitch covering theorem. It seems to me that this is very important and is certainly one of the most profound theorems I have ever seen. Differentiation theory is carried out for general Radon measures using this covering theorem. This is important because these kinds of measures are encountered in probability. Lebesgue measure is a special case. Abstract theory is presented later and includes the standard theorems on representation, Banach spaces, and so forth. Also included is a treatment of the Kolmogorov extension theorem. This major result is being neglected but, if I understand the situation correctly, it is the foundation for modern probability theory. It belongs in a course on analysis. The Bochner integral is also commonly neglected so I have included a careful treatment of this important topic. Some of us do our research in the context of spaces of Bochner integrable functions involving various function spaces.

There is an introduction to probability and stochastic processes at the end. I have included it because I encountered much of it in my old age and thought it was marvelous mathematics. I was not raised on it and this likely shows. However, it may be that someone can benefit from my efforts to understand this material. I have a hard time with it. There is more in my Topics in analysis book, but that is mostly pretty unorganized because I was gathering it from many different sources for our seminar. I am trying to present a more coherent presentation in this book. This is a very big topic and I must pick what I have found

most interesting, likely offending others who know better than I do what is important.

Chapter 1

Review of Some Linear Algebra

This material can be referred to as needed. It is here in order to make the book self contained.

1.1 The Matrix of a Linear Map

Recall the definition of a linear map. First of all, these need to be defined on a linear space and have values in a linear space.

Definition 1.1.1 Let $T : V \rightarrow W$ be a function. Here V and W are linear spaces. Then $T \in \mathcal{L}(V, W)$ and is a linear map means that for α, β scalars and v_1, v_2 vectors,

$$T(\alpha v_1 + \beta v_2) = \alpha T v_1 + \beta T v_2$$

Also recall from linear algebra that if you have $T \in \mathcal{L}(\mathbb{F}^n, \mathbb{F}^m)$ it can always be understood in terms of a matrix. That is, there exists an $m \times n$ matrix A such that for all $x \in \mathbb{F}^n$,

$$Ax = Tx$$

Recall that, from the way we multiply matrices,

$$A = (Te_1 \quad \cdots \quad Te_n)$$

That is, the i^{th} column is just Te_i .

1.2 Block Multiplication of Matrices

Consider the following problem

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix} \begin{pmatrix} E & F \\ G & H \end{pmatrix}.$$

You know how to do this. You get

$$\begin{pmatrix} AE + BG & AF + BH \\ CE + DG & CF + DH \end{pmatrix}.$$

Now what if instead of numbers, the entries, A, B, C, D, E, F, G are matrices of a size such that the multiplications and additions needed in the above formula all make sense. Would the formula be true in this case?

Suppose A is a matrix of the form

$$A = \begin{pmatrix} A_{11} & \cdots & A_{1m} \\ \vdots & \ddots & \vdots \\ A_{r1} & \cdots & A_{rm} \end{pmatrix} \tag{1.1}$$

where A_{ij} is a $s_i \times p_j$ matrix where s_i is constant for $j = 1, \dots, m$ for each $i = 1, \dots, r$. Such a matrix is called a **block matrix**, also a **partitioned matrix**. How do you get the block

A_{ij} ? Here is how for A an $m \times n$ matrix:

$$\overbrace{\begin{pmatrix} \mathbf{0} & I_{s_i \times s_i} & \mathbf{0} \end{pmatrix}}^{s_i \times m} A \overbrace{\begin{pmatrix} \mathbf{0} \\ I_{p_j \times p_j} \\ \mathbf{0} \end{pmatrix}}^{n \times p_j}. \quad (1.2)$$

In the block column matrix on the right, you need to have $c_j - 1$ rows of zeros above the small $p_j \times p_j$ identity matrix where the columns of A involved in A_{ij} are $c_j, \dots, c_j + p_j - 1$ and in the block row matrix on the left, you need to have $r_i - 1$ columns of zeros to the left of the $s_i \times s_i$ identity matrix where the rows of A involved in A_{ij} are $r_i, \dots, r_i + s_i$. An important observation to make is that the matrix on the right specifies columns to use in the block and the one on the left specifies the rows. Thus the block A_{ij} , in this case, is a matrix of size $s_i \times p_j$. There is no overlap between the blocks of A . Thus the identity $n \times n$ identity matrix corresponding to multiplication on the right of A is of the form

$$\begin{pmatrix} I_{p_1 \times p_1} & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & I_{p_m \times p_m} \end{pmatrix},$$

where these little identity matrices don't overlap. A similar conclusion follows from consideration of the matrices $I_{s_i \times s_i}$. Note that in (1.2), the matrix on the right is a block column matrix for the above block diagonal matrix, and the matrix on the left in (1.2) is a block row matrix taken from a similar block diagonal matrix consisting of the $I_{s_i \times s_i}$.

Next consider the question of multiplication of two block matrices. Let B be a block matrix of the form

$$\begin{pmatrix} B_{11} & \cdots & B_{1p} \\ \vdots & \ddots & \vdots \\ B_{r1} & \cdots & B_{rp} \end{pmatrix} \quad (1.3)$$

and A is a block matrix of the form

$$\begin{pmatrix} A_{11} & \cdots & A_{1m} \\ \vdots & \ddots & \vdots \\ A_{p1} & \cdots & A_{pm} \end{pmatrix} \quad (1.4)$$

such that for all i, j , it makes sense to multiply $B_{is}A_{sj}$ for all $s \in \{1, \dots, p\}$. (That is the two matrices B_{is} and A_{sj} are conformable.) and that for fixed ij , it follows that $B_{is}A_{sj}$ is the same size for each s so that it makes sense to write $\sum_s B_{is}A_{sj}$.

The following theorem says essentially that when you take the product of two matrices, you can partition both matrices, formally multiply the blocks to get another block matrix and this one will be BA partitioned. Before presenting this theorem, here is a simple lemma which is really a special case of the theorem.

Lemma 1.2.1 *Consider the following product.*

$$\begin{pmatrix} \mathbf{0} \\ I \\ \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{0} & I & \mathbf{0} \end{pmatrix}$$

where the first is $n \times r$ and the second is $r \times n$. The small identity matrix I is an $r \times r$ matrix and there are l zero rows above I and l zero columns to the left of I in the right matrix. Then the product of these matrices is a block matrix of the form

$$\begin{pmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & I & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix}.$$

Proof: From the definition of matrix multiplication, the product is

$$\left(\begin{pmatrix} \mathbf{0} \\ I \\ \mathbf{0} \end{pmatrix} \mathbf{0} \quad \cdots \quad \begin{pmatrix} \mathbf{0} \\ I \\ \mathbf{0} \end{pmatrix} e_1 \quad \cdots \quad \begin{pmatrix} \mathbf{0} \\ I \\ \mathbf{0} \end{pmatrix} e_r \quad \cdots \quad \begin{pmatrix} \mathbf{0} \\ I \\ \mathbf{0} \end{pmatrix} \mathbf{0} \right)$$

which yields the claimed result. In the formula e_j refers to the column vector of length r which has a 1 in the j^{th} position. This proves the lemma. ■

Theorem 1.2.2 Let B be a $q \times p$ block matrix as in (1.3) and let A be a $p \times n$ block matrix as in (1.4) such that B_{is} is conformable with A_{sj} and each product, $B_{is}A_{sj}$ for $s = 1, \dots, p$ is of the same size, so that they can be added. Then BA can be obtained as a block matrix such that the ij^{th} block is of the form

$$\sum_s B_{is}A_{sj}. \quad (1.5)$$

Proof: From (1.2)

$$B_{is}A_{sj} = \begin{pmatrix} \mathbf{0} & I_{r_i \times r_i} & \mathbf{0} \end{pmatrix} B \begin{pmatrix} \mathbf{0} \\ I_{p_s \times p_s} \\ \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{0} & I_{p_s \times p_s} & \mathbf{0} \end{pmatrix} A \begin{pmatrix} \mathbf{0} \\ I_{q_j \times q_j} \\ \mathbf{0} \end{pmatrix}$$

where here it is assumed B_{is} is $r_i \times p_s$ and A_{sj} is $p_s \times q_j$. The product involves the s^{th} block in the i^{th} row of blocks for B and the s^{th} block in the j^{th} column of A . Thus there are the same number of rows above the $I_{p_s \times p_s}$ as there are columns to the left of $I_{p_s \times p_s}$ in those two inside matrices. Then from Lemma 1.2.1

$$\begin{pmatrix} \mathbf{0} \\ I_{p_s \times p_s} \\ \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{0} & I_{p_s \times p_s} & \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & I_{p_s \times p_s} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix}.$$

Since the blocks of small identity matrices do not overlap,

$$\sum_s \begin{pmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & I_{p_s \times p_s} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix} = \begin{pmatrix} I_{p_1 \times p_1} & & 0 \\ & \ddots & \\ 0 & & I_{p_p \times p_p} \end{pmatrix} = I,$$

and so,

$$\sum_s B_{is}A_{sj} = \sum_s \begin{pmatrix} \mathbf{0} & I_{r_i \times r_i} & \mathbf{0} \end{pmatrix} B \begin{pmatrix} \mathbf{0} \\ I_{p_s \times p_s} \\ \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{0} & I_{p_s \times p_s} & \mathbf{0} \end{pmatrix} A \begin{pmatrix} \mathbf{0} \\ I_{q_j \times q_j} \\ \mathbf{0} \end{pmatrix}$$

$$\begin{aligned}
&= \begin{pmatrix} \mathbf{0} & I_{r_i \times r_i} & \mathbf{0} \end{pmatrix} B \sum_s \begin{pmatrix} \mathbf{0} \\ I_{p_s \times p_s} \\ \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{0} & I_{p_s \times p_s} & \mathbf{0} \end{pmatrix} A \begin{pmatrix} \mathbf{0} \\ I_{q_j \times q_j} \\ \mathbf{0} \end{pmatrix} \\
&= \begin{pmatrix} \mathbf{0} & I_{r_i \times r_i} & \mathbf{0} \end{pmatrix} B I A \begin{pmatrix} \mathbf{0} \\ I_{q_j \times q_j} \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{0} & I_{r_i \times r_i} & \mathbf{0} \end{pmatrix} B A \begin{pmatrix} \mathbf{0} \\ I_{q_j \times q_j} \\ \mathbf{0} \end{pmatrix}
\end{aligned}$$

which equals the ij^{th} block of BA . Hence the ij^{th} block of BA equals the formal multiplication according to matrix multiplication,

$$\sum_s B_{is} A_{sj}.$$

This proves the theorem. ■

Example 1.2.3 Multiply the following pair of partitioned matrices using the above theorem by multiplying the blocks as described above and then in the conventional manner.

$$\left(\begin{array}{cc|c} 1 & 2 & 3 \\ \hline -1 & 2 & 3 \\ 3 & -2 & 1 \end{array} \right) \left(\begin{array}{cc|c} 1 & -1 & 2 \\ \hline 2 & 3 & 0 \\ -2 & 2 & 1 \end{array} \right)$$

Doing it in terms of the blocks, this yields, after the indicated multiplications of the blocks,

$$\left(\begin{array}{cc|c} 5 + (-6) & & \\ \hline \begin{pmatrix} 3 \\ -1 \end{pmatrix} + \begin{pmatrix} 3 \\ 1 \end{pmatrix} (-2) & & \end{array} \right) \left(\begin{array}{cc|c} \begin{pmatrix} 5 & 2 \end{pmatrix} + 3 \begin{pmatrix} 2 & 1 \end{pmatrix} & & \\ \hline \begin{pmatrix} 7 & -2 \\ -9 & 6 \end{pmatrix} + \begin{pmatrix} 6 & 3 \\ 2 & 1 \end{pmatrix} & & \end{array} \right)$$

This is

$$\left(\begin{array}{cc|c} -1 & & \begin{pmatrix} 11 & 5 \end{pmatrix} \\ \hline \begin{pmatrix} -3 \\ -3 \end{pmatrix} & & \begin{pmatrix} 13 & 1 \\ -7 & 7 \end{pmatrix} \end{array} \right)$$

Multiplying it out the usual way, you have

$$\begin{pmatrix} 1 & 2 & 3 \\ -1 & 2 & 3 \\ 3 & -2 & 1 \end{pmatrix} \begin{pmatrix} 1 & -1 & 2 \\ 2 & 3 & 0 \\ -2 & 2 & 1 \end{pmatrix} = \begin{pmatrix} -1 & 11 & 5 \\ -3 & 13 & 1 \\ -3 & -7 & 7 \end{pmatrix}$$

you see this is the same thing without the partition lines.

1.3 Schur's Theorem

For some reason, not understood by me, Schur's theorem is often neglected in beginning linear algebra. This is too bad because it is one of the best theorems in linear algebra. Here $|\cdot|$ denotes the usual norm in \mathbb{C}^n given by

$$|x|^2 \equiv \sum_{j=1}^n |x_j|^2$$

Definition 1.3.1 A complex $n \times n$ matrix U is said to be unitary if $U^*U = I$. Here U^* is the transpose of the conjugate of U . The matrix is unitary if and only if its columns form an orthonormal set in \mathbb{C}^n . This follows from the way we multiply matrices in which the ij^{th} entry of U^*U is obtained by taking the conjugate of the i^{th} row of U times the j^{th} column of U .

Theorem 1.3.2 (Schur) Let A be a complex $n \times n$ matrix. Then there exists a unitary matrix U such that

$$U^*AU = T, \quad (1.6)$$

where T is an upper triangular matrix having the eigenvalues of A on the main diagonal, listed with multiplicity¹.

Proof: The theorem is clearly true if A is a 1×1 matrix. Just let $U = 1$, the 1×1 matrix which has entry 1. Suppose it is true for $(n-1) \times (n-1)$ matrices and let A be an $n \times n$ matrix. Then let v_1 be a unit eigenvector for A . Then there exists λ_1 such that

$$Av_1 = \lambda_1 v_1, \quad |v_1| = 1.$$

Extend $\{v_1\}$ to a basis and then use the Gram - Schmidt process to obtain

$$\{v_1, \dots, v_n\}$$

an orthonormal basis of \mathbb{C}^n . Let U_0 be a matrix whose i^{th} column is v_i . Then from the definition of a unitary matrix Definition 1.3.1, it follows that U_0 is unitary. Consider $U_0^*AU_0$.

$$U_0^*AU_0 = \begin{pmatrix} v_1^* \\ \vdots \\ v_n^* \end{pmatrix} \begin{pmatrix} Av_1 & \cdots & Av_n \end{pmatrix} = \begin{pmatrix} v_1^* \\ \vdots \\ v_n^* \end{pmatrix} \begin{pmatrix} \lambda_1 v_1 & \cdots & Av_n \end{pmatrix}$$

Thus $U_0^*AU_0$ is of the form

$$\begin{pmatrix} \lambda_1 & a \\ \mathbf{0} & A_1 \end{pmatrix}$$

where A_1 is an $(n-1) \times (n-1)$ matrix. Now by induction, there exists an $(n-1) \times (n-1)$ unitary matrix \tilde{U}_1 such that

$$\tilde{U}_1^* A_1 \tilde{U}_1 = T_{n-1},$$

an upper triangular matrix. Consider

$$U_1 \equiv \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \tilde{U}_1 \end{pmatrix}.$$

An application of block multiplication shows that U_1 is a unitary matrix and also that

$$U_1^* U_0^* A U_0 U_1 = \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \tilde{U}_1^* \end{pmatrix} \begin{pmatrix} \lambda_1 & * \\ \mathbf{0} & A_1 \end{pmatrix} \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \tilde{U}_1 \end{pmatrix} = \begin{pmatrix} \lambda_1 & * \\ \mathbf{0} & T_{n-1} \end{pmatrix} = T$$

where T is upper triangular. Then let $U = U_0 U_1$. Since $(U_0 U_1)^* = U_1^* U_0^*$, it follows that A is similar to T and that $U_0 U_1$ is unitary. Hence A and T have the same characteristic

¹ 'Listed with multiplicity' means that the diagonal entries are repeated according to their multiplicity as roots of the characteristic equation.

polynomials, and since the eigenvalues of T are the diagonal entries listed with multiplicity, this proves the theorem. ■

The same argument yields the following corollary in the case where A has real entries. The only difference is the use of the real inner product instead of the complex inner product.

Corollary 1.3.3 *Let A be a real $n \times n$ matrix which has only real eigenvalues. Then there exists a real orthogonal matrix Q such that*

$$Q^T A Q = T$$

where T is an upper triangular matrix having the eigenvalues of A on the main diagonal, listed with multiplicity.

Proof: This follows by observing that if all eigenvalues are real, then corresponding to each real eigenvalue, there exists a real eigenvector. Thus the argument of the above theorem applies with the real inner product in \mathbb{R}^n . ■

1.4 Hermitian and Symmetric Matrices

A complex $n \times n$ matrix A with $A^* = A$ is said to be **Hermitian**. A real $n \times n$ matrix A with $A^T = A$ is said to be **symmetric**. In either case, note that for $\langle \cdot, \cdot \rangle$ the inner product in \mathbb{C}^n ,

$$\langle A\mathbf{u}, \mathbf{v} \rangle = (A\mathbf{u})^T \bar{\mathbf{v}} = \mathbf{u}^T A^T \bar{\mathbf{v}} = \mathbf{u}^T \bar{A} \bar{\mathbf{v}} = \langle \mathbf{u}, A\mathbf{v} \rangle.$$

Thus, as a numerical example, the matrix

$$\begin{pmatrix} 1 & 1-i \\ 1+i & 2 \end{pmatrix}$$

is Hermitian, while

$$\begin{pmatrix} 1 & -1 & -2 \\ -1 & 2 & 4 \\ -2 & 4 & 3 \end{pmatrix}$$

is symmetric. Hermitian matrices are named in honor of the French mathematician Charles Hermite (1822–1901).

With Schur's theorem, the theorem on diagonalization of a Hermitian matrix follows.

Theorem 1.4.1 *Let A be Hermitian. Then the eigenvalues of A are all real, and there exists a unitary matrix U such that*

$$U^* A U = D,$$

a diagonal matrix whose diagonal entries are the eigenvalues of A listed with multiplicity. In case A is symmetric, U may be taken to be an orthogonal matrix. The columns of U form an orthonormal basis of eigenvectors of A .

Proof: By Schur's theorem and the assumption that A is Hermitian, there exists a triangular matrix T , whose diagonal entries are the eigenvalues of A listed with multiplicity, and a unitary matrix U such that

$$T = U^* A U = U^* A^* U = (U^* A U)^* = T^*.$$

It follows from this that T is a diagonal matrix and has all real entries down the main diagonal. Hence the eigenvalues of A are real. If A is symmetric (real and Hermitian) it follows from Corollary 1.3.3 that U may be taken to be orthogonal (The columns are an orthonormal set in the inner product of \mathbb{R}^n).

That the columns of U form an orthonormal basis of eigenvectors of A , follows right away from the definition of matrix multiplication which implies that if \mathbf{u}_i is a column of U , then $A\mathbf{u}_i = \text{column } i \text{ of } (UD) = \lambda_i \mathbf{u}_i$. ■

1.5 The Right Polar Factorization

The right polar factorization involves writing a matrix as a product of two other matrices, one which preserves distances and the other which stretches and distorts. First here are some lemmas which review and add to many of the topics discussed so far about adjoints and orthonormal sets and such things. This is of fundamental significance in geometric measure theory and also in continuum mechanics. Not surprisingly the stress should depend on the part which stretches and distorts. See [23].

Lemma 1.5.1 *Let A be a Hermitian matrix such that all its eigenvalues are nonnegative. Then there exists a Hermitian matrix $A^{1/2}$ such that $A^{1/2}$ has all nonnegative eigenvalues and $(A^{1/2})^2 = A$.*

Proof: Since A is Hermitian, there exists a diagonal matrix D having all real nonnegative entries and a unitary matrix U such that $A = U^*DU$. This is from Theorem 1.4.1 above. Then denote by $D^{1/2}$ the matrix which is obtained by replacing each diagonal entry of D with its square root. Thus $D^{1/2}D^{1/2} = D$. Then define

$$A^{1/2} \equiv U^*D^{1/2}U.$$

Then

$$(A^{1/2})^2 = U^*D^{1/2}UU^*D^{1/2}U = U^*DU = A.$$

Since $D^{1/2}$ is real,

$$(U^*D^{1/2}U)^* = U^*(D^{1/2})^*(U^*)^* = U^*D^{1/2}U$$

so $A^{1/2}$ is Hermitian. ■

Next it is helpful to recall the Gram Schmidt algorithm and observe a certain property stated in the next lemma.

Lemma 1.5.2 *Suppose $\{\mathbf{w}_1, \dots, \mathbf{w}_r, \mathbf{v}_{r+1}, \dots, \mathbf{v}_p\}$ is a linearly independent set of vectors such that $\{\mathbf{w}_1, \dots, \mathbf{w}_r\}$ is an orthonormal set of vectors. Then when the Gram Schmidt process is applied to the vectors in the given order, it will not change any of the $\mathbf{w}_1, \dots, \mathbf{w}_r$.*

Proof: Let $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ be the orthonormal set delivered by the Gram Schmidt process. Then $\mathbf{u}_1 = \mathbf{w}_1$ because by definition, $\mathbf{u}_1 \equiv \mathbf{w}_1/|\mathbf{w}_1| = \mathbf{w}_1$. Now suppose $\mathbf{u}_j = \mathbf{w}_j$ for all $j \leq k \leq r$. Then if $k < r$, consider the definition of \mathbf{u}_{k+1} .

$$\mathbf{u}_{k+1} \equiv \frac{\mathbf{w}_{k+1} - \sum_{j=1}^{k+1} (\mathbf{w}_{k+1}, \mathbf{u}_j) \mathbf{u}_j}{\left| \mathbf{w}_{k+1} - \sum_{j=1}^{k+1} (\mathbf{w}_{k+1}, \mathbf{u}_j) \mathbf{u}_j \right|}$$

By induction, $u_j = w_j$ and so this reduces to $w_{k+1}/|w_{k+1}| = w_{k+1}$. ■

This lemma immediately implies the following lemma.

Lemma 1.5.3 *Let V be a subspace of dimension p and let $\{w_1, \dots, w_r\}$ be an orthonormal set of vectors in V . Then this orthonormal set of vectors may be extended to an orthonormal basis for V ,*

$$\{w_1, \dots, w_r, y_{r+1}, \dots, y_p\}$$

Proof: First extend the given linearly independent set $\{w_1, \dots, w_r\}$ to a basis for V and then apply the Gram Schmidt theorem to the resulting basis. Since $\{w_1, \dots, w_r\}$ is orthonormal it follows from Lemma 1.5.2 the result is of the desired form, an orthonormal basis extending $\{w_1, \dots, w_r\}$. ■

Here is another lemma about preserving distance.

Lemma 1.5.4 *Suppose R is an $m \times n$ matrix with $m \geq n$ and R preserves distances. Then $R^*R = I$. Also, if R takes an orthonormal basis to an orthonormal set, then R must preserve distances.*

Proof: Since R preserves distances, $|Rx| = |x|$ for every x . Therefore from the axioms of the dot product,

$$\begin{aligned} & |x|^2 + |y|^2 + (x, y) + (y, x) = |x + y|^2 = (R(x + y), R(x + y)) \\ &= (Rx, Rx) + (Ry, Ry) + (Rx, Ry) + (Ry, Rx) \\ &= |x|^2 + |y|^2 + (R^*Rx, y) + (y, R^*Rx) \end{aligned}$$

and so for all x, y ,

$$(R^*Rx - x, y) + (y, R^*Rx - x) = 0$$

Hence for all x, y ,

$$\operatorname{Re}(R^*Rx - x, y) = 0$$

Now for a x, y given, choose $\alpha \in \mathbb{C}$ such that

$$\alpha(R^*Rx - x, y) = |(R^*Rx - x, y)|$$

Then

$$0 = \operatorname{Re}(R^*Rx - x, \bar{\alpha}y) = \operatorname{Re} \alpha(R^*Rx - x, y) = |(R^*Rx - x, y)|$$

Thus $|(R^*Rx - x, y)| = 0$ for all x, y because the given x, y were arbitrary. Let $y = R^*Rx - x$ to conclude that for all x ,

$$R^*Rx - x = 0$$

which says $R^*R = I$ since x is arbitrary.

Consider the last claim. Let $R: \mathbb{F}^n \rightarrow \mathbb{F}^m$ such that $\{u_1, \dots, u_n\}$ is an orthonormal basis for \mathbb{F}^n and $\{Ru_1, \dots, Ru_n\}$ is also an orthonormal set, then

$$\left| R \left(\sum_i x_i u_i \right) \right|^2 = \left| \sum_i x_i Ru_i \right|^2 = \sum_i |x_i|^2 = \left| \sum_i x_i u_i \right|^2 \quad \blacksquare$$

With this preparation, here is the big theorem about the right polar factorization.

Theorem 1.5.5 *Let F be an $m \times n$ matrix where $m \geq n$. Then there exists a Hermitian $n \times n$ matrix U which has all nonnegative eigenvalues and an $m \times n$ matrix R which satisfies $R^*R = I$ such that $F = RU$.*

Proof: Consider F^*F . This is a Hermitian matrix because

$$(F^*F)^* = F^*(F^*)^* = F^*F$$

Also the eigenvalues of the $n \times n$ matrix F^*F are all nonnegative. This is because if \mathbf{x} is an eigenvalue,

$$\lambda(\mathbf{x}, \mathbf{x}) = (F^*F\mathbf{x}, \mathbf{x}) = (F\mathbf{x}, F\mathbf{x}) \geq 0.$$

Therefore, by Lemma 1.5.1, there exists an $n \times n$ Hermitian matrix U having all nonnegative eigenvalues such that

$$U^2 = F^*F.$$

Consider the subspace $U(\mathbb{F}^n)$. Let $\{U\mathbf{x}_1, \dots, U\mathbf{x}_r\}$ be an orthonormal basis for

$$U(\mathbb{F}^n) \subseteq \mathbb{F}^n.$$

Note that $U(\mathbb{F}^n)$ might not be all of \mathbb{F}^n . Using Lemma 1.5.3, extend to an orthonormal basis for all of \mathbb{F}^n ,

$$\{U\mathbf{x}_1, \dots, U\mathbf{x}_r, \mathbf{y}_{r+1}, \dots, \mathbf{y}_n\}.$$

Next observe that $\{F\mathbf{x}_1, \dots, F\mathbf{x}_r\}$ is also an orthonormal set of vectors in \mathbb{F}^m . This is because

$$\begin{aligned} (F\mathbf{x}_k, F\mathbf{x}_j) &= (F^*F\mathbf{x}_k, \mathbf{x}_j) = (U^2\mathbf{x}_k, \mathbf{x}_j) \\ &= (U\mathbf{x}_k, U^*\mathbf{x}_j) = (U\mathbf{x}_k, U\mathbf{x}_j) = \delta_{jk} \end{aligned}$$

Therefore, from Lemma 1.5.3 again, this orthonormal set of vectors can be extended to an orthonormal basis for \mathbb{F}^m ,

$$\{F\mathbf{x}_1, \dots, F\mathbf{x}_r, \mathbf{z}_{r+1}, \dots, \mathbf{z}_m\}$$

Thus there are at least as many \mathbf{z}_k as there are \mathbf{y}_j . Now for $\mathbf{x} \in \mathbb{F}^n$, since

$$\{U\mathbf{x}_1, \dots, U\mathbf{x}_r, \mathbf{y}_{r+1}, \dots, \mathbf{y}_n\}$$

is an orthonormal basis for \mathbb{F}^n , there exist unique scalars,

$$c_1 \dots, c_r, d_{r+1}, \dots, d_n$$

such that

$$\mathbf{x} = \sum_{k=1}^r c_k U\mathbf{x}_k + \sum_{k=r+1}^n d_k \mathbf{y}_k$$

Define

$$R\mathbf{x} \equiv \sum_{k=1}^r c_k F\mathbf{x}_k + \sum_{k=r+1}^n d_k \mathbf{z}_k \quad (1.7)$$

Then also there exist scalars b_k such that

$$U\mathbf{x} = \sum_{k=1}^r b_k U\mathbf{x}_k$$

and so from 1.7,

$$RU \mathbf{x} = \sum_{k=1}^r b_k F \mathbf{x}_k = F \left(\sum_{k=1}^r b_k \mathbf{x}_k \right)$$

Is $F \left(\sum_{k=1}^r b_k \mathbf{x}_k \right) = F(\mathbf{x})$?

$$\begin{aligned} & \left(F \left(\sum_{k=1}^r b_k \mathbf{x}_k \right) - F(\mathbf{x}), F \left(\sum_{k=1}^r b_k \mathbf{x}_k \right) - F(\mathbf{x}) \right) \\ &= \left((F^* F) \left(\sum_{k=1}^r b_k \mathbf{x}_k - \mathbf{x} \right), \left(\sum_{k=1}^r b_k \mathbf{x}_k - \mathbf{x} \right) \right) \\ &= \left(U^2 \left(\sum_{k=1}^r b_k \mathbf{x}_k - \mathbf{x} \right), \left(\sum_{k=1}^r b_k \mathbf{x}_k - \mathbf{x} \right) \right) \\ &= \left(U \left(\sum_{k=1}^r b_k \mathbf{x}_k - \mathbf{x} \right), U \left(\sum_{k=1}^r b_k \mathbf{x}_k - \mathbf{x} \right) \right) \\ &= \left(\sum_{k=1}^r b_k U \mathbf{x}_k - U \mathbf{x}, \sum_{k=1}^r b_k U \mathbf{x}_k - U \mathbf{x} \right) = 0 \end{aligned}$$

Therefore, $F \left(\sum_{k=1}^r b_k \mathbf{x}_k \right) = F(\mathbf{x})$ and this shows $RU \mathbf{x} = F \mathbf{x}$. From 1.7 it follows that R maps an orthonormal set to an orthonormal set and so R preserves distances. Therefore, by Lemma 1.5.4 $R^* R = I$. ■

1.6 Elementary matrices

The elementary matrices result from doing a row operation to the identity matrix.

As before, everything will apply to matrices having coefficients in an arbitrary field of scalars, although we will mainly feature the real numbers in the examples.

Definition 1.6.1 *The row operations consist of the following*

1. Switch two rows.
2. Multiply a row by a nonzero number.
3. Replace a row by the same row added to a multiple of another row.

We refer to these as the row operations of type 1, 2, and 3 respectively.

The elementary matrices are given in the following definition.

Definition 1.6.2 *The elementary matrices consist of those matrices which result by applying a row operation to an identity matrix. Those which involve switching rows of the identity are called permutation matrices. More generally, a permutation matrix is a matrix which comes by permuting the rows of the identity matrix, not just switching two rows.*

As an example of why these elementary matrices are interesting, consider the following. Letting \mathbf{r}_i be the row vector of all zeros except for a 1 in the i^{th} slot,

$$\begin{pmatrix} \mathbf{r}_2 \\ \mathbf{r}_1 \\ \mathbf{r}_3 \end{pmatrix} \begin{pmatrix} a & b & c & d \\ x & y & z & w \\ f & g & h & i \end{pmatrix} = \begin{pmatrix} x & y & z & w \\ a & b & c & d \\ f & g & h & i \end{pmatrix}.$$

A 3×4 matrix was multiplied on the left by an elementary matrix which was obtained from row operation 1 applied to switching the first two rows of the identity matrix. This resulted in applying the operation 1 to the given matrix. This is what happens in general.

Now consider what these elementary matrices look like. They are obtained from switching a couple of rows of the identity matrix. First P_{ij} , which involves switching row i and row j of the identity where Let $i < j$. Then, as above, Then, as above, $P^{ij} =$

$$\begin{pmatrix} \mathbf{r}_1 \\ \vdots \\ \mathbf{r}_j \\ \vdots \\ \mathbf{r}_i \\ \vdots \\ \mathbf{r}_n \end{pmatrix}$$

where

$$\mathbf{r}_j = (0 \cdots 1 \cdots 0)$$

with the 1 in the j^{th} position from the left.

For P^{ij} this matrix which involves switching the i and j rows of the identity. Now consider what this does to a column vector.

$$\begin{pmatrix} \mathbf{r}_1 \\ \vdots \\ \mathbf{r}_j \\ \vdots \\ \mathbf{r}_i \\ \vdots \\ \mathbf{r}_n \end{pmatrix} \begin{pmatrix} v_1 \\ \vdots \\ v_i \\ \vdots \\ v_j \\ \vdots \\ v_n \end{pmatrix} = \begin{pmatrix} v_1 \\ \vdots \\ v_j \\ \vdots \\ v_i \\ \vdots \\ v_n \end{pmatrix}.$$

Now we try multiplication of a matrix on the left by this elementary matrix P^{ij} . Thus,

$$P^{ij} \begin{pmatrix} a_{11} & a_{12} & \cdots & \cdots & \cdots & \cdots & a_{1p} \\ \vdots & \vdots & & & & & \vdots \\ a_{i1} & a_{i2} & \cdots & \cdots & \cdots & \cdots & a_{ip} \\ \vdots & \vdots & & & & & \vdots \\ a_{j1} & a_{j2} & \cdots & \cdots & \cdots & \cdots & a_{jp} \\ \vdots & \vdots & & & & & \vdots \\ a_{n1} & a_{n2} & \cdots & \cdots & \cdots & \cdots & a_{np} \end{pmatrix}.$$

has the indicated columns listed in order:

$$\begin{pmatrix} P^{ij} \begin{pmatrix} a_{11} \\ \vdots \\ a_{i1} \\ \vdots \\ a_{j1} \\ \vdots \\ a_{n1} \end{pmatrix} & P^{ij} \begin{pmatrix} a_{12} \\ \vdots \\ a_{i2} \\ \vdots \\ a_{j2} \\ \vdots \\ a_{n2} \end{pmatrix} & \cdots & P^{ij} \begin{pmatrix} a_{1p} \\ \vdots \\ a_{ip} \\ \vdots \\ a_{jp} \\ \vdots \\ a_{np} \end{pmatrix} \end{pmatrix}$$

$$= \begin{pmatrix} \begin{pmatrix} a_{11} \\ \vdots \\ a_{j1} \\ \vdots \\ a_{i1} \\ \vdots \\ a_{n1} \end{pmatrix} & \begin{pmatrix} a_{12} \\ \vdots \\ a_{j2} \\ \vdots \\ a_{i2} \\ \vdots \\ a_{n2} \end{pmatrix} & \cdots & \begin{pmatrix} a_{1p} \\ \vdots \\ a_{jp} \\ \vdots \\ a_{ip} \\ \vdots \\ a_{np} \end{pmatrix} \end{pmatrix}$$

and so the resulting matrix is

$$= \begin{pmatrix} a_{11} & a_{12} & \cdots & \cdots & \cdots & \cdots & a_{1p} \\ \vdots & \vdots & & & & & \vdots \\ a_{j1} & a_{j2} & \cdots & \cdots & \cdots & \cdots & a_{jp} \\ \vdots & \vdots & & & & & \vdots \\ a_{i1} & a_{i2} & \cdots & \cdots & \cdots & \cdots & a_{ip} \\ \vdots & \vdots & & & & & \vdots \\ a_{n1} & a_{n2} & \cdots & \cdots & \cdots & \cdots & a_{np} \end{pmatrix}.$$

This has established the following lemma.

Lemma 1.6.3 *Let P^{ij} denote the elementary matrix which involves switching the i^{th} and the j^{th} rows of I . Then if P^{ij} , A are conformable, we have*

$$P^{ij}A = B$$

where B is obtained from A by switching the i^{th} and the j^{th} rows.

Next consider the row operation which involves multiplying the i^{th} row by a nonzero constant, c . We write

$$I = \begin{pmatrix} \mathbf{r}_1 \\ \mathbf{r}_2 \\ \vdots \\ \mathbf{r}_n \end{pmatrix}$$

where

$$\mathbf{r}_j = (0 \cdots 1 \cdots 0)$$

with the 1 in the j^{th} position from the left. The elementary matrix which results from applying this operation to the i^{th} row of the identity matrix is of the form

$$E(c, i) = \begin{pmatrix} \mathbf{r}_1 \\ \vdots \\ c\mathbf{r}_i \\ \vdots \\ \mathbf{r}_n \end{pmatrix}.$$

Now consider what this does to a column vector.

$$\begin{pmatrix} \mathbf{r}_1 \\ \vdots \\ c\mathbf{r}_i \\ \vdots \\ \mathbf{r}_n \end{pmatrix} \begin{pmatrix} v_1 \\ \vdots \\ v_i \\ \vdots \\ v_n \end{pmatrix} = \begin{pmatrix} v_1 \\ \vdots \\ cv_i \\ \vdots \\ v_n \end{pmatrix}.$$

Denote by $E(c, i)$ this elementary matrix which multiplies the i^{th} row of the identity by the nonzero constant, c . Then from what was just discussed and the way matrices are multiplied,

$$E(c, i) \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1p} \\ \vdots & \vdots & & \vdots \\ a_{i1} & a_{i2} & \cdots & a_{ip} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{np} \end{pmatrix}$$

equals a matrix having the columns indicated below.

$$= \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1p} \\ \vdots & \vdots & & \vdots \\ ca_{i1} & ca_{i2} & \cdots & ca_{ip} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{np} \end{pmatrix}.$$

This proves the following lemma.

Lemma 1.6.4 *Let $E(c, i)$ denote the elementary matrix corresponding to the row operation in which the i^{th} row is multiplied by the nonzero constant c . Thus $E(c, i)$ involves multiplying the i^{th} row of the identity matrix by c . Then*

$$E(c, i)A = B$$

where B is obtained from A by multiplying the i^{th} row of A by c .

Finally consider the third of these row operations. Letting \mathbf{r}_j be the j^{th} row of the identity matrix, denote by $E(c \times i + j)$ the elementary matrix obtained from the identity

matrix by replacing \mathbf{r}_j with $\mathbf{r}_j + c\mathbf{r}_i$. In case $i < j$ this will be of the form

$$P^{ij} = \begin{pmatrix} \mathbf{r}_1 \\ \vdots \\ \mathbf{r}_i \\ \vdots \\ c\mathbf{r}_i + \mathbf{r}_j \\ \vdots \\ \mathbf{r}_n \end{pmatrix}.$$

Consider what this does to a column vector.

$$\begin{pmatrix} \mathbf{r}_1 \\ \vdots \\ \mathbf{r}_i \\ \vdots \\ c\mathbf{r}_i + \mathbf{r}_j \\ \vdots \\ \mathbf{r}_n \end{pmatrix} \begin{pmatrix} v_1 \\ \vdots \\ v_i \\ \vdots \\ v_j \\ \vdots \\ v_n \end{pmatrix} = \begin{pmatrix} v_1 \\ \vdots \\ v_i \\ \vdots \\ cv_i + v_j \\ \vdots \\ v_n \end{pmatrix}.$$

From this and the way matrices are multiplied,

$$E(c \times i + j) \begin{pmatrix} a_{11} & a_{12} & \cdots & \cdots & \cdots & \cdots & a_{1p} \\ \vdots & \vdots & & & & & \vdots \\ a_{i1} & a_{i2} & \cdots & \cdots & \cdots & \cdots & a_{ip} \\ \vdots & \vdots & & & & & \vdots \\ a_{j2} & a_{j2} & \cdots & \cdots & \cdots & \cdots & a_{jp} \\ \vdots & \vdots & & & & & \vdots \\ a_{n1} & a_{n2} & \cdots & \cdots & \cdots & \cdots & a_{np} \end{pmatrix}$$

equals a matrix having the indicated columns listed in order.

$$\left(E(c \times i + j) \begin{pmatrix} a_{11} \\ \vdots \\ a_{i1} \\ \vdots \\ a_{j2} \\ \vdots \\ a_{n1} \end{pmatrix}, E(c \times i + j) \begin{pmatrix} a_{12} \\ \vdots \\ a_{i2} \\ \vdots \\ a_{j2} \\ \vdots \\ a_{n2} \end{pmatrix}, \cdots E(c \times i + j) \begin{pmatrix} a_{1p} \\ \vdots \\ a_{ip} \\ \vdots \\ a_{jp} \\ \vdots \\ a_{np} \end{pmatrix} \right)$$

$$= \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1p} \\ \vdots & \vdots & & \vdots \\ a_{i1} & a_{i2} & \cdots & a_{ip} \\ \vdots & \vdots & & \vdots \\ a_{j2} + ca_{i1} & a_{j2} + ca_{i2} & \cdots & a_{jp} + ca_{ip} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{np} \end{pmatrix}.$$

The case where $i > j$ is similar. This proves the following lemma in which, as above, the i^{th} row of the identity is \mathbf{r}_i .

Lemma 1.6.5 *Let $E(c \times i + j)$ denote the elementary matrix obtained from I by replacing the j^{th} row of the identity \mathbf{r}_j with $c\mathbf{r}_i + \mathbf{r}_j$. Letting the k^{th} row of A be \mathbf{a}_k ,*

$$E(c \times i + j)A = B$$

where B has the same rows as A except the j^{th} row of B is $c\mathbf{a}_i + \mathbf{a}_j$.

The above lemmas are summarized in the following theorem.

Theorem 1.6.6 *To perform any of the three row operations on a matrix A it suffices to do the row operation on the identity matrix, obtaining an elementary matrix E , and then take the product, EA . In addition to this, the following identities hold for the elementary matrices described above.*

$$E(c \times i + j)E(-c \times i + j) = E(-c \times i + j)E(c \times i + j) = I. \quad (1.8)$$

$$E(c, i)E(c^{-1}, i) = E(c^{-1}, i)E(c, i) = I. \quad (1.9)$$

$$P^{ij}P^{ij} = I. \quad (1.10)$$

Proof: Consider (1.8). Starting with I and taking $-c$ times the i^{th} row added to the j^{th} yields $E(-c \times i + j)$ which differs from I only in the j^{th} row. Now multiplying on the left by $E(c \times i + j)$ takes c times the i^{th} row and adds to the j^{th} thus restoring the j^{th} row to its original state. Thus $E(c \times i + j)E(-c \times i + j) = I$. Similarly $E(-c \times i + j)E(c \times i + j) = I$. The reasoning is similar for (1.9) and (1.10). ■

Each of these elementary matrices has a significant geometric significance. The effect of doing $E(\frac{1}{2} \times 3 + 1)$ shears the box in one direction. Of course there would be corresponding shears in the other directions also. Note that this does not change the volume. You should think about the geometric effect of the other elementary matrices on a box.

Definition 1.6.7 *For an $n \times n$ matrix A , an $n \times n$ matrix B which has the property that $AB = BA = I$ is denoted by A^{-1} . Such a matrix is called an **inverse**. When A has an inverse, it is called **invertible**.*

The following lemma says that if a matrix acts like an inverse, then it is **the** inverse. Also, the product of invertible matrices is invertible.

Lemma 1.6.8 *If B, C are both inverses of A , then $B = C$. That is, there exists at most one inverse of a matrix. If A_1, \dots, A_m are each invertible $m \times m$ matrices, then the product $A_1 A_2 \cdots A_m$ is also invertible and*

$$(A_1 A_2 \cdots A_m)^{-1} = A_m^{-1} A_{m-1}^{-1} \cdots A_1^{-1}.$$

Proof. From the definition and associative law of matrix multiplication,

$$B = BI = B(AC) = (BA)C = IC = C.$$

This proves the uniqueness of the inverse.

Next suppose A, B are invertible. Then

$$AB(B^{-1}A^{-1}) = A(BB^{-1})A^{-1} = AIA^{-1} = AA^{-1} = I$$

and also

$$(B^{-1}A^{-1})AB = B^{-1}(A^{-1}A)B = B^{-1}IB = B^{-1}B = I.$$

It follows from Definition 1.6.7 that AB has an inverse and it is $B^{-1}A^{-1}$. Thus the case of $m = 1, 2$ in the claim of the lemma is true. Suppose this claim is true for k . Then

$$A_1 A_2 \cdots A_k A_{k+1} = (A_1 A_2 \cdots A_k) A_{k+1}.$$

By induction, the two matrices $(A_1 A_2 \cdots A_k), A_{k+1}$ are both invertible and

$$(A_1 A_2 \cdots A_k)^{-1} = A_k^{-1} \cdots A_2^{-1} A_1^{-1}.$$

By the case of the product of two invertible matrices shown above,

$$((A_1 A_2 \cdots A_k) A_{k+1})^{-1} = A_{k+1}^{-1} (A_1 A_2 \cdots A_k)^{-1} = A_{k+1}^{-1} A_k^{-1} \cdots A_2^{-1} A_1^{-1}.$$

This proves the lemma. ■

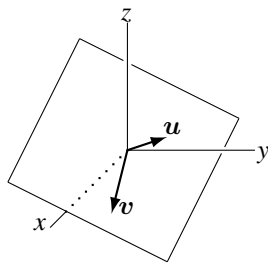
We will discuss methods for finding the inverse later. For now, observe that Theorem 1.6.6 says that elementary matrices are invertible and that the inverse of such a matrix is also an elementary matrix. The major conclusion of the above Lemma and Theorem is the following lemma about linear relationships.

Definition 1.6.9 *Let v_1, \dots, v_k, u be vectors. Then u is said to be a **linear combination** of the vectors $\{v_1, \dots, v_k\}$ if there exist scalars c_1, \dots, c_k such that*

$$u = \sum_{i=1}^k c_i v_i.$$

*We also say that when the above holds for some scalars c_1, \dots, c_k , there exists a **linear relationship** between the vector u and the vectors $\{v_1, \dots, v_k\}$.*

We will discuss this more later, but the following picture illustrates the geometric significance of the vectors which have a linear relationship with two vectors u, v pointing in different directions.



The following lemma states that linear relationships between columns in a matrix are preserved by row operations. This simple lemma is the main result in understanding all the major questions related to the row reduced echelon form as well as many other topics.

Lemma 1.6.10 *Let A and B be two $m \times n$ matrices and suppose B results from a row operation applied to A . Then the k^{th} column of B is a linear combination of the i_1, \dots, i_r columns of B if and only if the k^{th} column of A is a linear combination of the i_1, \dots, i_r columns of A . Furthermore, the scalars in the linear combinations are the same. (The linear relationship between the k^{th} column of A and the i_1, \dots, i_r columns of A is the same as the linear relationship between the k^{th} column of B and the i_1, \dots, i_r columns of B .)*

Proof: Let A be the following matrix in which the \mathbf{a}_k are the columns

$$\begin{pmatrix} \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_n \end{pmatrix}$$

and let B be the following matrix in which the columns are given by the \mathbf{b}_k

$$\begin{pmatrix} \mathbf{b}_1 & \mathbf{b}_2 & \cdots & \mathbf{b}_n \end{pmatrix}.$$

Then by Theorem 1.6.6 on Page 29, $\mathbf{b}_k = E\mathbf{a}_k$ where E is an elementary matrix. Suppose then that one of the columns of A is a linear combination of some other columns of A . Say

$$\mathbf{a}_k = c_1\mathbf{a}_{i_1} + \cdots + c_r\mathbf{a}_{i_r}.$$

Then multiplying by E ,

$$\mathbf{b}_k = E\mathbf{a}_k = c_1E\mathbf{a}_{i_1} + \cdots + c_rE\mathbf{a}_{i_r} = c_1\mathbf{b}_{i_1} + \cdots + c_r\mathbf{b}_{i_r}.$$

This proves the lemma. ■

Example 1.6.11 Find linear relationships between the columns of the matrix

$$A = \begin{pmatrix} 1 & 3 & 11 & 10 & 36 \\ 1 & 2 & 8 & 9 & 23 \\ 1 & 1 & 5 & 8 & 10 \end{pmatrix}.$$

It is not clear what the relationships are, so we do row operations to this matrix. Lemma 1.6.10 says that all the linear relationships between columns are preserved, so the idea is to do row operations until a matrix results which has the property that the linear relationships are obvious. First take -1 times the top row and add to the two bottom rows. This yields

$$\begin{pmatrix} 1 & 3 & 11 & 10 & 36 \\ 0 & -1 & -3 & -1 & -13 \\ 0 & -2 & -6 & -2 & -26 \end{pmatrix}$$

Next take -2 times the middle row and add to the bottom row followed by multiplying the middle row by -1 :

$$\begin{pmatrix} 1 & 3 & 11 & 10 & 36 \\ 0 & 1 & 3 & 1 & 13 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}.$$

Next take -3 times the middle row added to the top:

$$\begin{pmatrix} 1 & 0 & 2 & 7 & -3 \\ 0 & 1 & 3 & 1 & 13 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}. \quad (1.11)$$

At this point it is clear that the last column is -3 times the first column added to 13 times the second. By Lemma 1.6.10, the same is true of the corresponding columns in the original matrix A . As a check,

$$-3 \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} + 13 \begin{pmatrix} 3 \\ 2 \\ 1 \end{pmatrix} = \begin{pmatrix} 36 \\ 23 \\ 10 \end{pmatrix}.$$

You should notice that other linear relationships are also easily seen from (1.11). For example the fourth column is 7 times the first added to the second. This is obvious from (1.11) and Lemma 1.6.10 says the same relationship holds for A .

This is really just an extension of the technique for finding solutions to a linear system of equations. In solving a system of equations earlier, row operations were used to exhibit the last column of an augmented matrix as a linear combination of the preceding columns. The **row reduced echelon form** just extends this by making obvious the linear relationships between **every** column, not just the last, and those columns preceding it. The matrix in 1.11 is in row reduced echelon form. The row reduced echelon form is the topic of the next section.

1.7 The Row Reduced Echelon Form Of A Matrix

When you do row operations on a matrix, there is an ultimate conclusion. It is called the **row reduced echelon form**. We show here that every matrix has such a row reduced echelon form and that this row reduced echelon form is unique. The significance is that it becomes possible to use the definite article in referring to **the** row reduced echelon form. Hence important conclusions about the original matrix may be logically deduced from an examination of its unique row reduced echelon form. First we need the following definition.

Definition 1.7.1 Define special column vectors e_i as follows.

$$e_i = (0 \quad \cdots \quad 1 \quad \cdots \quad 0)^T.$$

Recall that T says to take the transpose. Thus e_i is the column vector which has all zero entries except for a 1 in the i^{th} position down from the top.

Now here is the description of the row reduced echelon form.

Definition 1.7.2 An $m \times n$ matrix is said to be in **row reduced echelon form** if, in viewing successive columns from left to right, the first nonzero column encountered is

e_1 and if, in viewing the columns of the matrix from left to right, you have encountered e_1, e_2, \dots, e_k , the next column is either e_{k+1} or this next column is a linear combination of the vectors, e_1, e_2, \dots, e_k .

Example 1.7.3 The following matrices are in row reduced echelon form.

$$\begin{pmatrix} 1 & 0 & 4 & 0 \\ 0 & 1 & 3 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}, \begin{pmatrix} 0 & 1 & 0 & 0 & 7 \\ 0 & 0 & 0 & 1 & 3 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}, \begin{pmatrix} 0 & 1 & 0 & 3 \\ 0 & 0 & 1 & -5 \\ 0 & 0 & 0 & 0 \end{pmatrix}.$$

Definition 1.7.4 Given a matrix A , row reduction produces one and only one row reduced matrix B with $A \sim B$. See Corollary 1.7.9. We call B **the** row reduced echelon form of A .

Theorem 1.7.5 Let A be an $m \times n$ matrix. Then A has a row reduced echelon form determined by a simple process.

Proof. Viewing the columns of A from left to right, take the first nonzero column. Pick a nonzero entry in this column and switch the row containing this entry with the top row of A . Now divide this new top row by the value of this nonzero entry to get a 1 in this position and then use row operations to make all entries below this equal to zero. Thus the first nonzero column is now e_1 . Denote the resulting matrix by A_1 . Consider the sub-matrix of A_1 to the right of this column and below the first row. Do exactly the same thing for this sub-matrix that was done for A . This time the e_1 will refer to F^{m-1} . Use the first 1 obtained by the above process which is in the top row of this sub-matrix and row operations, to produce a zero in place of every entry above it and below it. Call the resulting matrix A_2 . Thus A_2 satisfies the conditions of the above definition up to the column just encountered. Continue this way till every column has been dealt with and the result must be in row reduced echelon form. ■

Here is some terminology about pivot columns.

Definition 1.7.6 The first **pivot column** of A is the first nonzero column of A which becomes e_1 in the row reduced echelon form. The next pivot column is the first column after this which becomes e_2 in the row reduced echelon form. The third is the next column which becomes e_3 in the row reduced echelon form and so forth.

The algorithm just described for obtaining a row reduced echelon form shows that these columns are well defined, but we will deal with this issue more carefully in Corollary 1.7.9 where we show that every matrix corresponds to exactly one row reduced echelon form.

Definition 1.7.7 Two matrices A, B are said to be **row equivalent** if B can be obtained from A by a sequence of row operations. When A is row equivalent to B , we write $A \sim B$.

Proposition 1.7.8 In the notation of Definition 1.7.7. $A \sim A$. If $A \sim B$, then $B \sim A$. If $A \sim B$ and $B \sim C$, then $A \sim C$.

Proof. That $A \sim A$ is obvious. Consider the second claim. By Theorem 1.6.6, there exist elementary matrices E_1, E_2, \dots, E_m such that

$$B = E_1 E_2 \cdots E_m A.$$

It follows from Lemma 1.6.8 that $(E_1 E_2 \cdots E_m)^{-1}$ exists and equals the product of the inverses of these matrices in the reverse order. Thus

$$\begin{aligned} E_m^{-1} E_{m-1}^{-1} \cdots E_1^{-1} B &= (E_1 E_2 \cdots E_m)^{-1} B \\ &= (E_1 E_2 \cdots E_m)^{-1} (E_1 E_2 \cdots E_m) A = A. \end{aligned}$$

By Theorem 1.6.6, each E_k^{-1} is an elementary matrix. By Theorem 1.6.6 again, the above shows that A results from a sequence of row operations applied to B . The last claim is left for an exercise. This proves the proposition. ■

There are three choices for row operations at each step in Theorem 1.7.5. A natural question is whether the same row reduced echelon matrix always results in the end from following any sequence of row operations.

We have already made use of the following observation in finding a linear relationship between the columns of the matrix A , but here it is stated more formally.

$$\begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} = x_1 e_1 + \cdots + x_n e_n,$$

so to say two column vectors are equal, is to say the column vectors are the same linear combination of the special vectors e_j .

Corollary 1.7.9 *The row reduced echelon form is unique. That is if B, C are two matrices in row reduced echelon form and both are obtained from A by a sequence of row operations, then $B = C$.*

Proof. Suppose B and C are both row reduced echelon forms for the matrix A . It follows that B and C have zero columns in the same positions because row operations do not affect zero columns. By Proposition 1.7.8, B and C are row equivalent. In reading from left to right in B , suppose e_1, \dots, e_r occur first in positions i_1, \dots, i_r respectively. Then from the description of the row reduced echelon form, each of these columns of B , in positions i_1, \dots, i_r , is not a linear combination of the preceding columns. Since C is row equivalent to B , it follows from Lemma 1.6.10, that each column of C in positions i_1, \dots, i_r is not a linear combination of the preceding columns of C . By the description of the row reduced echelon form, e_1, \dots, e_r occur first in C , in positions i_1, \dots, i_r respectively. Therefore, both B and C have the sequence e_1, e_2, \dots, e_r occurring first (reading from left to right) in the positions, i_1, i_2, \dots, i_r . Since these matrices are row equivalent, it follows from Lemma 1.6.10, that the columns between the i_k and i_{k+1} position in the two matrices are linear combinations involving the same scalars, of the columns in the i_1, \dots, i_k position. Similarly, the columns after the i_r position are linear combinations of the columns in the i_1, \dots, i_r positions involving the same scalars in both matrices. This is equivalent to the assertion that each of these columns is identical in B and C . ■

Now with the above corollary, here is a very fundamental observation. The number of nonzero rows in the row reduced echelon form is the same as the number of pivot columns.

Namely, this number is r in both cases where e_1, \dots, e_r are the pivot columns in the row reduced echelon form. This number r is called the **rank** of the matrix. This is discussed more later, but first here are some other applications.

Consider a matrix which looks like this: (More columns than rows.)



Corollary 1.7.10 Suppose A is an $m \times n$ matrix and that $m < n$. That is, the number of rows is less than the number of columns. Then one of the columns of A is a linear combination of the preceding columns of A . Also, there exists $\mathbf{x} \in F^n$ such that $\mathbf{x} \neq \mathbf{0}$ and $A\mathbf{x} = \mathbf{0}$.

Proof: Since $m < n$, not all the columns of A can be pivot columns. In reading from left to right, pick the first one which is not a pivot column. Then from the description of the row reduced echelon form, this column is a linear combination of the preceding columns. Say

$$\mathbf{a}_j = x_1 \mathbf{a}_1 + \cdots + x_{j-1} \mathbf{a}_{j-1}.$$

Therefore, from the way we multiply a matrix times a vector,

$$A \begin{pmatrix} x_1 \\ \vdots \\ x_{j-1} \\ -1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = (\mathbf{a}_1 \cdots \mathbf{a}_{j-1} \mathbf{a}_j \cdots \mathbf{a}_n) \begin{pmatrix} x_1 \\ \vdots \\ x_{j-1} \\ -1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \mathbf{0}. \blacksquare$$

1.8 Finding the Inverse of a Matrix

Recall that the inverse of an $n \times n$ matrix A is a matrix B such that

$$AB = BA = I$$

where I is the identity matrix. It was shown that an elementary matrix is invertible and that its inverse is also an elementary matrix. Also the product of invertible matrices is invertible and its inverse is the product of the inverses in the reverse order. In this section, we consider the problem of finding an inverse for a given $n \times n$ matrix.

Example 1.8.1 Let $A = \begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix}$. Show that $\begin{pmatrix} 2 & -1 \\ -1 & 1 \end{pmatrix}$ is the inverse of A .

To check this, multiply

$$\begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} 2 & -1 \\ -1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix},$$

and

$$\begin{pmatrix} 2 & -1 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix},$$

showing that this matrix is indeed the inverse of A .

In the last example, how would you find A^{-1} ? You wish to find a matrix $\begin{pmatrix} x & z \\ y & w \end{pmatrix}$ such that

$$\begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} x & z \\ y & w \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}.$$

This requires the solution of the systems of equations,

$$x + y = 1, x + 2y = 0$$

and

$$z + w = 0, z + 2w = 1.$$

Writing the augmented matrix for these two systems gives

$$\left(\begin{array}{cc|c} 1 & 1 & 1 \\ 1 & 2 & 0 \end{array} \right) \quad (1.12)$$

for the first system and

$$\left(\begin{array}{cc|c} 1 & 1 & 0 \\ 1 & 2 & 1 \end{array} \right) \quad (1.13)$$

for the second. Let's solve the first system. Take (-1) times the first row and add to the second to get

$$\left(\begin{array}{cc|c} 1 & 1 & 1 \\ 0 & 1 & -1 \end{array} \right)$$

Now take (-1) times the second row and add to the first to get

$$\left(\begin{array}{cc|c} 1 & 0 & 2 \\ 0 & 1 & -1 \end{array} \right).$$

Putting in the variables, this says $x = 2$ and $y = -1$.

Now solve the second system, (1.13) to find z and w . Take (-1) times the first row and add to the second to get

$$\left(\begin{array}{cc|c} 1 & 1 & 0 \\ 0 & 1 & 1 \end{array} \right).$$

Now take (-1) times the second row and add to the first to get

$$\left(\begin{array}{cc|c} 1 & 0 & -1 \\ 0 & 1 & 1 \end{array} \right).$$

Putting in the variables, this says $z = -1$ and $w = 1$. Therefore, the inverse is

$$\begin{pmatrix} 2 & -1 \\ -1 & 1 \end{pmatrix}.$$

Didn't the above seem rather repetitive? Exactly the same row operations were used in both systems. In each case, the end result was something of the form $(I|v)$ where I is the

identity and v gave a column of the inverse. In the above $\begin{pmatrix} x \\ y \end{pmatrix}$, the first column of the inverse was obtained first and then the second column $\begin{pmatrix} z \\ w \end{pmatrix}$.

To simplify this procedure, you could have written

$$\left(\begin{array}{cc|cc} 1 & 1 & 1 & 0 \\ 1 & 2 & 0 & 1 \end{array} \right)$$

and row reduced till you obtained

$$\left(\begin{array}{cc|cc} 1 & 0 & 2 & -1 \\ 0 & 1 & -1 & 1 \end{array} \right).$$

Then you could have read off the inverse as the 2×2 matrix on the right side. You should be able to see that it is valid by adapting the argument used in the simple case above.

This is the reason for the following simple procedure for finding the inverse of a matrix. This procedure is called the **Gauss-Jordan procedure**.

Procedure 1.8.2 Suppose A is an $n \times n$ matrix. To find A^{-1} if it exists, form the augmented $n \times 2n$ matrix

$$(A|I)$$

and then if possible, do row operations until you obtain an $n \times 2n$ matrix of the form

$$(I|B). \quad (1.14)$$

When this has been done, $B = A^{-1}$. If it is impossible to row reduce to a matrix of the form $(I|B)$, then A has no inverse.

The procedure just described along with the preceding explanation shows that this procedure actually yields a **right inverse**. This is a matrix B such that $AB = I$. We will show in Theorem 1.8.4 that this right inverse is really **the** inverse. This is a stronger result than that of Lemma 1.6.8 about the uniqueness of the inverse. For now, here is an example.

Example 1.8.3 Let $A = \begin{pmatrix} 1 & 2 & 2 \\ 1 & 0 & 2 \\ 3 & 1 & -1 \end{pmatrix}$. Find A^{-1} if it exists.

Set up the augmented matrix $(A|I)$:

$$\left(\begin{array}{ccc|ccc} 1 & 2 & 2 & 1 & 0 & 0 \\ 1 & 0 & 2 & 0 & 1 & 0 \\ 3 & 1 & -1 & 0 & 0 & 1 \end{array} \right)$$

Next take (-1) times the first row and add to the second followed by (-3) times the first row added to the last. This yields

$$\left(\begin{array}{ccc|ccc} 1 & 2 & 2 & 1 & 0 & 0 \\ 0 & -2 & 0 & -1 & 1 & 0 \\ 0 & -5 & -7 & -3 & 0 & 1 \end{array} \right).$$

Then take 5 times the second row and add to -2 times the last row.

$$\left(\begin{array}{ccc|ccc} 1 & 2 & 2 & 1 & 0 & 0 \\ 0 & -10 & 0 & -5 & 5 & 0 \\ 0 & 0 & 14 & 1 & 5 & -2 \end{array} \right)$$

Next take the last row and add to (-7) times the top row. This yields

$$\left(\begin{array}{ccc|ccc} -7 & -14 & 0 & -6 & 5 & -2 \\ 0 & -10 & 0 & -5 & 5 & 0 \\ 0 & 0 & 14 & 1 & 5 & -2 \end{array} \right).$$

Now take $(-7/5)$ times the second row and add to the top.

$$\left(\begin{array}{ccc|ccc} -7 & 0 & 0 & 1 & -2 & -2 \\ 0 & -10 & 0 & -5 & 5 & 0 \\ 0 & 0 & 14 & 1 & 5 & -2 \end{array} \right).$$

Finally divide the top row by -7 , the second row by -10 and the bottom row by 14 , which yields

$$\left(\begin{array}{ccc|ccc} 1 & 0 & 0 & -\frac{1}{7} & \frac{2}{7} & \frac{2}{7} \\ 0 & 1 & 0 & \frac{1}{2} & -\frac{1}{2} & 0 \\ 0 & 0 & 1 & \frac{1}{14} & \frac{5}{14} & -\frac{1}{7} \end{array} \right).$$

Therefore, the inverse is

$$\left(\begin{array}{ccc} -\frac{1}{7} & \frac{2}{7} & \frac{2}{7} \\ \frac{1}{2} & -\frac{1}{2} & 0 \\ \frac{1}{14} & \frac{5}{14} & -\frac{1}{7} \end{array} \right).$$

What you have really found in the above algorithm is a **right inverse**. Is this right inverse matrix, which we have called the inverse, really **the** inverse, the matrix which when multiplied on both sides gives the identity?

Theorem 1.8.4 Suppose A, B are $n \times n$ matrices and $AB = I$. Then it follows that $BA = I$ also, and so $B = A^{-1}$. For $n \times n$ matrices, the left inverse, right inverse and inverse are all the same thing.

Proof. If $AB = I$ for A, B $n \times n$ matrices, is $BA = I$? If $AB = I$, there exists a unique solution x to the equation

$$Bx = y$$

for any choice of y . In fact,

$$x = A(Bx) = Ay.$$

This means the row reduced echelon form of B must be I . Thus every column is a pivot column. Otherwise, there exists a free variable and the solution, if it exists, would not be

unique, contrary to what was just shown must happen if $AB = I$. It follows that a right inverse B^{-1} for B exists. The above procedure yields

$$(B \ I) \rightarrow (I \ B^{-1}).$$

Now multiply both sides of the equation $AB = I$ on the right by B^{-1} . Then

$$A = A(BB^{-1}) = (AB)B^{-1} = B^{-1}.$$

Thus A is the right inverse of B , and so $BA = I$. This shows that if $AB = I$, then $BA = I$ also. Exchanging roles of A and B , we see that if $BA = I$, then $AB = I$. This proves the theorem. ■

This has shown that in the context of $n \times n$ matrices, right inverses, left inverses and inverses are all the same and this matrix is called A^{-1} .

The following corollary is also of interest.

Corollary 1.8.5 *An $n \times n$ matrix A has an inverse if and only if the row reduced echelon form of A is I .*

Proof. First suppose the row reduced echelon form of A is I . Then Procedure 1.8.2 yields a right inverse for A . By Theorem 1.8.4 this is **the** inverse. Next suppose A has an inverse. Then there exists a unique solution \mathbf{x} to the equation $A\mathbf{x} = \mathbf{y}$, given by $\mathbf{x} = A^{-1}\mathbf{y}$. It follows that in the augmented matrix $(A|0)$ there are no free variables, and so every column to the left of the zero column is a pivot column. Therefore, the row reduced echelon form of A is I . ■

1.9 The Mathematical Theory of Determinants

It is easiest to give a different definition of the determinant which is clearly well defined and then prove the earlier one in terms of Laplace expansion. Let (i_1, \dots, i_n) be an ordered list of numbers from $\{1, \dots, n\}$. This means the order is important so $(1, 2, 3)$ and $(2, 1, 3)$ are different. There will be some repetition between this section and the earlier section on determinants. The main purpose is to give all the missing proofs. Two books which give a good introduction to determinants are Apostol [1] and Rudin [49]. A recent book which also has a good introduction is Baker [4].

1.9.1 The Function sgn

The following Lemma will be essential in the definition of the determinant.

Lemma 1.9.1 *There exists a function, sgn_n which maps each ordered list of numbers from $\{1, \dots, n\}$ to one of the three numbers, 0, 1, or -1 which also has the following properties.*

$$\text{sgn}_n(1, \dots, n) = 1 \tag{1.15}$$

$$\text{sgn}_n(i_1, \dots, p, \dots, q, \dots, i_n) = -\text{sgn}_n(i_1, \dots, q, \dots, p, \dots, i_n) \tag{1.16}$$

In words, the second property states that if two of the numbers are switched, the value of the function is multiplied by -1 . Also, in the case where $n > 1$ and $\{i_1, \dots, i_n\} = \{1, \dots, n\}$ so that every number from $\{1, \dots, n\}$ appears in the ordered list, (i_1, \dots, i_n) ,

$$\text{sgn}_n(i_1, \dots, i_{\theta-1}, n, i_{\theta+1}, \dots, i_n) \equiv$$

$$(-1)^{n-\theta} \operatorname{sgn}_{n-1}(i_1, \dots, i_{\theta-1}, i_{\theta+1}, \dots, i_n) \quad (1.17)$$

where $n = i_\theta$ in the ordered list, (i_1, \dots, i_n) .

Proof: Define $\operatorname{sign}(x) = 1$ if $x > 0$, -1 if $x < 0$ and 0 if $x = 0$. If $n = 1$, there is only one list and it is just the number 1. Thus one can define $\operatorname{sgn}_1(1) \equiv 1$. For the general case where $n > 1$, simply define

$$\operatorname{sgn}_n(i_1, \dots, i_n) \equiv \operatorname{sign} \left(\prod_{r < s} (i_s - i_r) \right)$$

This delivers either $-1, 1$, or 0 by definition. What about the other claims? Suppose you switch i_p with i_q where $p < q$ so two numbers in the ordered list (i_1, \dots, i_n) are switched. Denote the new ordered list of numbers as (j_1, \dots, j_n) . Thus $j_p = i_q$ and $j_q = i_p$ and if $r \notin \{p, q\}$, $j_r = i_r$. See the following illustration

i_1	i_2	\dots	i_p	\dots	i_q	\dots	i_n
1	2	\dots	p	\dots	q	\dots	n
i_1	i_2	\dots	i_q	\dots	i_p	\dots	i_n
1	2	\dots	q	\dots	p	\dots	n
j_1	j_2	\dots	j_p	\dots	j_q	\dots	j_n
1	2	\dots	p	\dots	q	\dots	n

Then

$$\begin{aligned} \operatorname{sgn}_n(j_1, \dots, j_n) &\equiv \operatorname{sign} \left(\prod_{r < s} (j_s - j_r) \right) \\ &= \operatorname{sign} \left(\overset{\text{both } p, q}{(i_p - i_q)} \prod_{p < j < q} \overset{\text{one of } p, q}{(i_j - i_q)} \prod_{p < j < q} \overset{\text{neither } p \text{ nor } q}{(i_p - i_j)} \prod_{r < s, r, s \notin \{p, q\}} (i_s - i_r) \right) \end{aligned}$$

The last product consists of the product of terms which were in $\prod_{r < s} (i_s - i_r)$ while the two products in the middle both introduce $q - p - 1$ minus signs. Thus their product is positive. The first factor is of opposite sign to the $i_q - i_p$ which occurred in $\operatorname{sgn}_n(i_1, \dots, i_n)$. Therefore, this switch introduced a minus sign and

$$\operatorname{sgn}_n(j_1, \dots, j_n) = -\operatorname{sgn}_n(i_1, \dots, i_n)$$

Now consider the last claim. In computing $\operatorname{sgn}_n(i_1, \dots, i_{\theta-1}, n, i_{\theta+1}, \dots, i_n)$ there will be the product of $n - \theta$ negative terms

$$(i_{\theta+1} - n) \cdots (i_n - n)$$

and the other terms in the product for computing $\operatorname{sgn}_n(i_1, \dots, i_{\theta-1}, n, i_{\theta+1}, \dots, i_n)$ are those which are required to compute $\operatorname{sgn}_{n-1}(i_1, \dots, i_{\theta-1}, i_{\theta+1}, \dots, i_n)$ multiplied by terms of the form $(n - i_j)$ which are nonnegative. It follows that

$$\operatorname{sgn}_n(i_1, \dots, i_{\theta-1}, n, i_{\theta+1}, \dots, i_n) = (-1)^{n-\theta} \operatorname{sgn}_{n-1}(i_1, \dots, i_{\theta-1}, i_{\theta+1}, \dots, i_n)$$

It is obvious that if there are repeats in the list the function gives 0. ■

Lemma 1.9.2 *Every ordered list of distinct numbers from $\{1, 2, \dots, n\}$ can be obtained from every other ordered list of distinct numbers by a finite number of switches. Also, sgn_n is unique.*

Proof: This is obvious if $n = 1$ or 2 . Suppose then that it is true for sets of $n - 1$ elements. Take two ordered lists of numbers, P_1, P_2 . Make one switch in both to place n at the end. Call the result P_1^n and P_2^n . Then using induction, there are finitely many switches in P_1^n so that it will coincide with P_2^n . Now switch the n in what results to where it was in P_2 .

To see sgn_n is unique, if there exist two functions, f and g both satisfying 1.15 and 1.16, you could start with $f(1, \dots, n) = g(1, \dots, n) = 1$ and applying the same sequence of switches, eventually arrive at $f(i_1, \dots, i_n) = g(i_1, \dots, i_n)$. If any numbers are repeated, then 1.16 gives both functions are equal to zero for that ordered list. ■

Definition 1.9.3 *Given an ordered list of distinct numbers from $\{1, 2, \dots, n\}$, say*

$$(i_1, \dots, i_n),$$

this ordered list is called a permutation. The symbol for all such permutations is S_n . The number $\text{sgn}_n(i_1, \dots, i_n)$ is called the sign of the permutation.

A permutation can also be considered as a function from the set

$$\{1, 2, \dots, n\} \text{ to } \{1, 2, \dots, n\}$$

as follows. Let $f(k) = i_k$. Permutations are of fundamental importance in certain areas of math. For example, it was by considering permutations that Galois was able to give a criterion for solution of polynomial equations by radicals, but this is a different direction than what is being attempted here.

In what follows sgn will often be used rather than sgn_n because the context supplies the appropriate n .

1.9.2 The Definition of the Determinant

Definition 1.9.4 *Let f be a real valued function which has the set of ordered lists of numbers from $\{1, \dots, n\}$ as its domain. Define*

$$\sum_{(k_1, \dots, k_n)} f(k_1 \dots k_n)$$

to be the sum of all the $f(k_1 \dots k_n)$ for all possible choices of ordered lists (k_1, \dots, k_n) of numbers of $\{1, \dots, n\}$. For example,

$$\sum_{(k_1, k_2)} f(k_1, k_2) = f(1, 2) + f(2, 1) + f(1, 1) + f(2, 2).$$

Definition 1.9.5 *Let $(a_{ij}) = A$ denote an $n \times n$ matrix. The determinant of A , denoted by $\det(A)$ is defined by*

$$\det(A) \equiv \sum_{(k_1, \dots, k_n)} \text{sgn}(k_1, \dots, k_n) a_{1k_1} \dots a_{nk_n}$$

where the sum is taken over all ordered lists of numbers from $\{1, \dots, n\}$. Note it suffices to take the sum over only those ordered lists in which there are no repeats because if there are, $\text{sgn}(k_1, \dots, k_n) = 0$ and so that term contributes 0 to the sum.

Let A be an $n \times n$ matrix $A = (a_{ij})$ and let (r_1, \dots, r_n) denote an ordered list of n numbers from $\{1, \dots, n\}$. Let $A(r_1, \dots, r_n)$ denote the matrix whose k^{th} row is the r_k row of the matrix A . Thus

$$\det(A(r_1, \dots, r_n)) = \sum_{(k_1, \dots, k_n)} \text{sgn}(k_1, \dots, k_n) a_{r_1 k_1} \cdots a_{r_n k_n} \quad (1.18)$$

and $A(1, \dots, n) = A$.

Proposition 1.9.6 *Let (r_1, \dots, r_n) be an ordered list of numbers from*

$$\{1, \dots, n\}$$

Then

$$\text{sgn}(r_1, \dots, r_n) \det(A) = \sum_{(k_1, \dots, k_n)} \text{sgn}(k_1, \dots, k_n) a_{r_1 k_1} \cdots a_{r_n k_n} \quad (1.19)$$

$$= \det(A(r_1, \dots, r_n)). \quad (1.20)$$

Proof: Let $(1, \dots, n) = (1, \dots, r, \dots, s, \dots, n)$ so $r < s$.

$$\det(A(1, \dots, r, \dots, s, \dots, n)) = \quad (1.21)$$

$$\sum_{(k_1, \dots, k_n)} \text{sgn}(k_1, \dots, k_r, \dots, k_s, \dots, k_n) a_{1k_1} \cdots a_{rk_r} \cdots a_{sk_s} \cdots a_{nk_n},$$

and renaming the variables, calling k_s, k_r and k_r, k_s , this equals

$$\begin{aligned} &= \sum_{(k_1, \dots, k_n)} \text{sgn}(k_1, \dots, k_s, \dots, k_r, \dots, k_n) a_{1k_1} \cdots a_{rk_s} \cdots a_{sk_r} \cdots a_{nk_n} \\ &= \sum_{(k_1, \dots, k_n)} -\text{sgn} \left(k_1, \dots, \overbrace{k_r, \dots, k_s}^{\text{These got switched}}, \dots, k_n \right) a_{1k_1} \cdots a_{sk_r} \cdots a_{rk_s} \cdots a_{nk_n} \\ &= -\det(A(1, \dots, s, \dots, r, \dots, n)). \end{aligned} \quad (1.22)$$

Consequently,

$$\det(A(1, \dots, s, \dots, r, \dots, n)) = -\det(A(1, \dots, r, \dots, s, \dots, n)) = -\det(A)$$

Now letting $A(1, \dots, s, \dots, r, \dots, n)$ play the role of A , and continuing in this way, switching pairs of numbers,

$$\det(A(r_1, \dots, r_n)) = (-1)^p \det(A)$$

where it took p switches to obtain (r_1, \dots, r_n) from $(1, \dots, n)$. By Lemma 1.9.1, this implies

$$\det(A(r_1, \dots, r_n)) = (-1)^p \det(A) = \text{sgn}(r_1, \dots, r_n) \det(A)$$

and proves the proposition in the case when there are no repeated numbers in the ordered list, (r_1, \dots, r_n) . However, if there is a repeat, say the r^{th} row equals the s^{th} row, then the reasoning of 1.21 -1.22 shows that $\det(A(r_1, \dots, r_n)) = 0$ and also $\text{sgn}(r_1, \dots, r_n) = 0$ so the formula holds in this case also. ■

Observation 1.9.7 *There are $n!$ ordered lists of distinct numbers from*

$$\{1, \dots, n\}$$

To see this, consider n slots placed in order. There are n choices for the first slot. For each of these choices, there are $n - 1$ choices for the second. Thus there are $n(n - 1)$ ways to fill the first two slots. Then for each of these ways there are $n - 2$ choices left for the third slot. Continuing this way, there are $n!$ ordered lists of distinct numbers from $\{1, \dots, n\}$ as stated in the observation.

1.9.3 A Symmetric Definition

With the above, it is possible to give a more symmetric description of the determinant from which it will follow that $\det(A) = \det(A^T)$.

Corollary 1.9.8 *The following formula for $\det(A)$ is valid.*

$$\det(A) = \frac{1}{n!} \cdot \sum_{(r_1, \dots, r_n)} \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(r_1, \dots, r_n) \operatorname{sgn}(k_1, \dots, k_n) a_{r_1 k_1} \cdots a_{r_n k_n}. \quad (1.23)$$

And also $\det(A^T) = \det(A)$ where A^T is the transpose of A . (Recall that for $A^T = (a_{ij}^T)$, $a_{ij}^T = a_{ji}$.)

Proof: From Proposition 1.9.6, if the r_i are distinct,

$$\det(A) = \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(r_1, \dots, r_n) \operatorname{sgn}(k_1, \dots, k_n) a_{r_1 k_1} \cdots a_{r_n k_n}.$$

Summing over all ordered lists, (r_1, \dots, r_n) where the r_i are distinct, (If the r_i are not distinct, $\operatorname{sgn}(r_1, \dots, r_n) = 0$ and so there is no contribution to the sum.)

$$n! \det(A) = \sum_{(r_1, \dots, r_n)} \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(r_1, \dots, r_n) \operatorname{sgn}(k_1, \dots, k_n) a_{r_1 k_1} \cdots a_{r_n k_n}.$$

This proves the corollary since the formula gives the same number for A as it does for A^T . ■

Corollary 1.9.9 *If two rows or two columns in an $n \times n$ matrix A , are switched, the determinant of the resulting matrix equals (-1) times the determinant of the original matrix. If A is an $n \times n$ matrix in which two rows are equal or two columns are equal then $\det(A) = 0$. Suppose the i^{th} row of A equals*

$$(xa_1 + yb_1, \dots, xa_n + yb_n)$$

Then

$$\det(A) = x \det(A_1) + y \det(A_2)$$

where the i^{th} row of A_1 is (a_1, \dots, a_n) and the i^{th} row of A_2 is (b_1, \dots, b_n) , all other rows of A_1 and A_2 coinciding with those of A . In other words, \det is a linear function of each row A . The same is true with the word “row” replaced with the word “column”.

Proof: By Proposition 1.9.6 when two rows are switched, the determinant of the resulting matrix is (-1) times the determinant of the original matrix. By Corollary 1.9.8 the same holds for columns because the columns of the matrix equal the rows of the transposed matrix. Thus if A_1 is the matrix obtained from A by switching two columns,

$$\det(A) = \det(A^T) = -\det(A_1^T) = -\det(A_1).$$

If A has two equal columns or two equal rows, then switching them results in the same matrix. Therefore, $\det(A) = -\det(A)$ and so $\det(A) = 0$.

It remains to verify the last assertion.

$$\begin{aligned} \det(A) &\equiv \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_n) a_{1k_1} \cdots (xa_{rk_i} + yb_{rk_i}) \cdots a_{nk_n} \\ &= x \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_n) a_{1k_1} \cdots a_{rk_i} \cdots a_{nk_n} \\ &\quad + y \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_n) a_{1k_1} \cdots b_{rk_i} \cdots a_{nk_n} \equiv x \det(A_1) + y \det(A_2). \end{aligned}$$

The same is true of columns because $\det(A^T) = \det(A)$ and the rows of A^T are the columns of A . ■

1.9.4 Basic Properties of the Determinant

Definition 1.9.10 A vector, \mathbf{w} , is a linear combination $\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$ if there exist scalars c_1, \dots, c_r such that $\mathbf{w} = \sum_{k=1}^r c_k \mathbf{v}_k$. This is the same as saying

$$\mathbf{w} \in \operatorname{span}(\mathbf{v}_1, \dots, \mathbf{v}_r).$$

The following corollary is also of great use.

Corollary 1.9.11 Suppose A is an $n \times n$ matrix and some column (row) is a linear combination of r other columns (rows). Then $\det(A) = 0$.

Proof: Let $A = (\mathbf{a}_1 \cdots \mathbf{a}_n)$ be the columns of A and suppose the condition that one column is a linear combination of r of the others is satisfied. Say $\mathbf{a}_i = \sum_{j \neq i} c_j \mathbf{a}_j$. Then by Corollary 1.9.9, $\det(A) =$

$$\det(\mathbf{a}_1 \cdots \sum_{j \neq i} c_j \mathbf{a}_j \cdots \mathbf{a}_n) = \sum_{j \neq i} c_j \det(\mathbf{a}_1 \cdots \mathbf{a}_j \cdots \mathbf{a}_n) = 0$$

because each of these determinants in the sum has two equal rows. ■

Recall the following definition of matrix multiplication.

Definition 1.9.12 If A and B are $n \times n$ matrices, $A = (a_{ij})$ and $B = (b_{ij})$, $AB = (c_{ij})$ where $c_{ij} \equiv \sum_{k=1}^n a_{ik} b_{kj}$.

One of the most important rules about determinants is that the determinant of a product equals the product of the determinants.

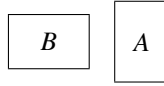
Theorem 1.9.13 *Let A and B be $n \times n$ matrices. Then*

$$\det(AB) = \det(A) \det(B).$$

Proof: Let c_{ij} be the ij^{th} entry of AB . Then by Proposition 1.9.6,

$$\begin{aligned} \det(AB) &= \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_n) c_{1k_1} \cdots c_{nk_n} \\ &= \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_n) \left(\sum_{r_1} a_{1r_1} b_{r_1 k_1} \right) \cdots \left(\sum_{r_n} a_{nr_n} b_{r_n k_n} \right) \\ &= \sum_{(r_1, \dots, r_n)} \sum_{(k_1, \dots, k_n)} \operatorname{sgn}(k_1, \dots, k_n) b_{r_1 k_1} \cdots b_{r_n k_n} (a_{1r_1} \cdots a_{nr_n}) \\ &= \sum_{(r_1, \dots, r_n)} \operatorname{sgn}(r_1 \cdots r_n) a_{1r_1} \cdots a_{nr_n} \det(B) = \det(A) \det(B). \blacksquare \end{aligned}$$

The Binet Cauchy formula is a generalization of the theorem which says the determinant of a product is the product of the determinants. The situation is illustrated in the following picture where A, B are matrices.



Theorem 1.9.14 *Let A be an $n \times m$ matrix with $n \geq m$ and let B be a $m \times n$ matrix. Also let A_i*

$$i = 1, \dots, C(n, m)$$

be the $m \times m$ submatrices of A which are obtained by deleting $n - m$ rows and let B_i be the $m \times m$ submatrices of B which are obtained by deleting corresponding $n - m$ columns. Then

$$\det(BA) = \sum_{k=1}^{C(n, m)} \det(B_k) \det(A_k)$$

Proof: This follows from a computation. By Corollary 1.9.8 on Page 43, $\det(BA) =$

$$\begin{aligned} &\frac{1}{m!} \sum_{(i_1 \cdots i_m)} \sum_{(j_1 \cdots j_m)} \operatorname{sgn}(i_1 \cdots i_m) \operatorname{sgn}(j_1 \cdots j_m) (BA)_{i_1 j_1} (BA)_{i_2 j_2} \cdots (BA)_{i_m j_m} \\ &= \frac{1}{m!} \sum_{(i_1 \cdots i_m)} \sum_{(j_1 \cdots j_m)} \operatorname{sgn}(i_1 \cdots i_m) \operatorname{sgn}(j_1 \cdots j_m) \cdot \\ &\quad \sum_{r_1=1}^n B_{i_1 r_1} A_{r_1 j_1} \sum_{r_2=1}^n B_{i_2 r_2} A_{r_2 j_2} \cdots \sum_{r_m=1}^n B_{i_m r_m} A_{r_m j_m} \end{aligned}$$

Now denote by I_k one of the subsets of $\{1, \dots, n\}$ which has m elements. Thus there are $C(n, m)$ of these.

$$\begin{aligned} &= \sum_{k=1}^{C(n, m)} \sum_{\{r_1, \dots, r_m\} = I_k} \frac{1}{m!} \sum_{(i_1 \cdots i_m)} \sum_{(j_1 \cdots j_m)} \operatorname{sgn}(i_1 \cdots i_m) \operatorname{sgn}(j_1 \cdots j_m) \cdot \\ &\quad B_{i_1 r_1} A_{r_1 j_1} B_{i_2 r_2} A_{r_2 j_2} \cdots B_{i_m r_m} A_{r_m j_m} \end{aligned}$$

$$\begin{aligned}
&= \sum_{k=1}^{C(n,m)} \sum_{\{r_1, \dots, r_m\}=I_k} \frac{1}{m!} \sum_{(i_1 \dots i_m)} \operatorname{sgn}(i_1 \dots i_m) B_{i_1 r_1} B_{i_2 r_2} \dots B_{i_m r_m} \cdot \\
&\quad \sum_{(j_1 \dots j_m)} \operatorname{sgn}(j_1 \dots j_m) A_{r_1 j_1} A_{r_2 j_2} \dots A_{r_m j_m} \\
&= \sum_{k=1}^{C(n,m)} \sum_{\{r_1, \dots, r_m\}=I_k} \frac{1}{m!} \operatorname{sgn}(r_1 \dots r_m)^2 \det(B_k) \det(A_k) = \sum_{k=1}^{C(n,m)} \det(B_k) \det(A_k)
\end{aligned}$$

since there are $m!$ ways of arranging the indices $\{r_1, \dots, r_m\}$. ■

1.9.5 Expansion Using Cofactors

Lemma 1.9.15 Suppose a matrix is of the form

$$M = \begin{pmatrix} A & * \\ \mathbf{0} & a \end{pmatrix} \text{ or } \begin{pmatrix} A & \mathbf{0} \\ * & a \end{pmatrix} \quad (1.24)$$

where a is a number and A is an $(n-1) \times (n-1)$ matrix and $*$ denotes either a column or a row having length $n-1$ and the $\mathbf{0}$ denotes either a column or a row of length $n-1$ consisting entirely of zeros. Then $\det(M) = a \det(A)$.

Proof: Denote M by (m_{ij}) . Thus in the first case, $m_{nn} = a$ and $m_{ni} = 0$ if $i \neq n$ while in the second case, $m_{nn} = a$ and $m_{in} = 0$ if $i \neq n$. From the definition of the determinant,

$$\det(M) \equiv \sum_{(k_1, \dots, k_n)} \operatorname{sgn}_n(k_1, \dots, k_n) m_{1k_1} \dots m_{nk_n}$$

Letting θ denote the position of n in the ordered list, (k_1, \dots, k_n) then using the earlier conventions used to prove Lemma 1.9.1, $\det(M)$ equals

$$\sum_{(k_1, \dots, k_n)} (-1)^{n-\theta} \operatorname{sgn}_{n-1} \left(k_1, \dots, k_{\theta-1}, k_{\theta+1}, \dots, k_n \right) m_{1k_1} \dots m_{nk_n}$$

Now suppose the second case. Then if $k_n \neq n$, the term involving m_{nk_n} in the above expression equals zero. Therefore, the only terms which survive are those for which $\theta = n$ or in other words, those for which $k_n = n$. Therefore, the above expression reduces to

$$a \sum_{(k_1, \dots, k_{n-1})} \operatorname{sgn}_{n-1}(k_1, \dots, k_{n-1}) m_{1k_1} \dots m_{(n-1)k_{n-1}} = a \det(A).$$

To get the assertion in the first case, use Corollary 1.9.8 to write

$$\det(M) = \det(M^T) = \det \left(\begin{pmatrix} A^T & \mathbf{0} \\ * & a \end{pmatrix} \right) = a \det(A^T) = a \det(A). \blacksquare$$

In terms of the theory of determinants, arguably the most important idea is that of Laplace expansion along a row or a column. This will follow from the above definition of a determinant.

Definition 1.9.16 Let $A = (a_{ij})$ be an $n \times n$ matrix. Then a new matrix called the *cofactor matrix* $\text{cof}(A)$ is defined by $\text{cof}(A) = (c_{ij})$ where to obtain c_{ij} delete the i^{th} row and the j^{th} column of A , take the determinant of the $(n-1) \times (n-1)$ matrix which results, (This is called the ij^{th} minor of A .) and then multiply this number by $(-1)^{i+j}$. To make the formulas easier to remember, $\text{cof}(A)_{ij}$ will denote the ij^{th} entry of the cofactor matrix.

The following is the main result. Earlier this was given as a definition and the outrageous totally unjustified assertion was made that the same number would be obtained by expanding the determinant along any row or column. The following theorem proves this assertion.

Theorem 1.9.17 Let A be an $n \times n$ matrix where $n \geq 2$. Then

$$\det(A) = \sum_{j=1}^n a_{ij} \text{cof}(A)_{ij} = \sum_{i=1}^n a_{ij} \text{cof}(A)_{ij}. \quad (1.25)$$

The first formula consists of expanding the determinant along the i^{th} row and the second expands the determinant along the j^{th} column.

Proof: Let (a_{i1}, \dots, a_{in}) be the i^{th} row of A . Let B_j be the matrix obtained from A by leaving every row the same except the i^{th} row which in B_j equals $(0, \dots, 0, a_{ij}, 0, \dots, 0)$. Then by Corollary 1.9.9,

$$\det(A) = \sum_{j=1}^n \det(B_j)$$

For example if

$$A = \begin{pmatrix} a & b & c \\ d & e & f \\ h & i & j \end{pmatrix}$$

and $i = 2$, then

$$B_1 = \begin{pmatrix} a & b & c \\ d & 0 & 0 \\ h & i & j \end{pmatrix}, B_2 = \begin{pmatrix} a & b & c \\ 0 & e & 0 \\ h & i & j \end{pmatrix}, B_3 = \begin{pmatrix} a & b & c \\ 0 & 0 & f \\ h & i & j \end{pmatrix}$$

Denote by A^{ij} the $(n-1) \times (n-1)$ matrix obtained by deleting the i^{th} row and the j^{th} column of A . Thus $\text{cof}(A)_{ij} \equiv (-1)^{i+j} \det(A^{ij})$. At this point, recall that from Proposition 1.9.6, when two rows or two columns in a matrix M , are switched, this results in multiplying the determinant of the old matrix by -1 to get the determinant of the new matrix. Therefore, by Lemma 1.9.15,

$$\begin{aligned} \det(B_j) &= (-1)^{n-j} (-1)^{n-i} \det \left(\begin{pmatrix} A^{ij} & * \\ \mathbf{0} & a_{ij} \end{pmatrix} \right) \\ &= (-1)^{i+j} \det \left(\begin{pmatrix} A^{ij} & * \\ \mathbf{0} & a_{ij} \end{pmatrix} \right) = a_{ij} \text{cof}(A)_{ij}. \end{aligned}$$

Therefore,

$$\det(A) = \sum_{j=1}^n a_{ij} \text{cof}(A)_{ij}$$

which is the formula for expanding $\det(A)$ along the i^{th} row. Also,

$$\det(A) = \det(A^T) = \sum_{j=1}^n a_{ij}^T \operatorname{cof}(A^T)_{ij} = \sum_{j=1}^n a_{ji} \operatorname{cof}(A)_{ji}$$

which is the formula for expanding $\det(A)$ along the i^{th} column. ■

1.9.6 A Formula for the Inverse

Note that this gives an easy way to write a formula for the inverse of an $n \times n$ matrix. Recall the definition of the inverse of a matrix in Definition 1.6.7 on Page 29.

Theorem 1.9.18 A^{-1} exists if and only if $\det(A) \neq 0$. If $\det(A) \neq 0$, then $A^{-1} = (a_{ij}^{-1})$ where

$$a_{ij}^{-1} = \det(A)^{-1} \operatorname{cof}(A)_{ji}$$

for $\operatorname{cof}(A)_{ij}$ the ij^{th} cofactor of A .

Proof: By Theorem 1.9.17 and letting $(a_{ir}) = A$, if $\det(A) \neq 0$,

$$\sum_{i=1}^n a_{ir} \operatorname{cof}(A)_{ir} \det(A)^{-1} = \det(A) \det(A)^{-1} = 1.$$

Now in the matrix A , replace the k^{th} column with the r^{th} column and then expand along the k^{th} column. This yields for $k \neq r$,

$$\sum_{i=1}^n a_{ir} \operatorname{cof}(A)_{ik} \det(A)^{-1} = 0$$

because there are two equal columns by Corollary 1.9.9. Summarizing,

$$\sum_{i=1}^n a_{ir} \operatorname{cof}(A)_{ik} \det(A)^{-1} = \delta_{rk}.$$

Using the other formula in Theorem 1.9.17, and similar reasoning,

$$\sum_{j=1}^n a_{rj} \operatorname{cof}(A)_{kj} \det(A)^{-1} = \delta_{rk}$$

This proves that if $\det(A) \neq 0$, then A^{-1} exists with $A^{-1} = (a_{ij}^{-1})$, where

$$a_{ij}^{-1} = \operatorname{cof}(A)_{ji} \det(A)^{-1}.$$

Now suppose A^{-1} exists. Then by Theorem 1.9.13,

$$1 = \det(I) = \det(AA^{-1}) = \det(A) \det(A^{-1})$$

so $\det(A) \neq 0$. ■

The next corollary points out that if an $n \times n$ matrix A has a right or a left inverse, then it has an inverse.

Corollary 1.9.19 *Let A be an $n \times n$ matrix and suppose there exists an $n \times n$ matrix B such that $BA = I$. Then A^{-1} exists and $A^{-1} = B$. Also, if there exists C an $n \times n$ matrix such that $AC = I$, then A^{-1} exists and $A^{-1} = C$.*

Proof: Since $BA = I$, Theorem 1.9.13 implies $\det B \det A = 1$ and so $\det A \neq 0$. Therefore from Theorem 1.9.18, A^{-1} exists. Therefore,

$$A^{-1} = (BA)A^{-1} = B(AA^{-1}) = BI = B.$$

The case where $CA = I$ is handled similarly. ■

The conclusion of this corollary is that left inverses, right inverses and inverses are all the same in the context of $n \times n$ matrices.

Theorem 1.9.18 says that to find the inverse, take the transpose of the cofactor matrix and divide by the determinant. The transpose of the cofactor matrix is called the adjugate or sometimes the classical adjoint of the matrix A . It is an abomination to call it the adjoint although you do sometimes see it referred to in this way. In words, A^{-1} is equal to one over the determinant of A times the adjugate matrix of A .

1.9.7 Cramer's Rule

In case you are solving a system of equations, $Ax = y$ for x , it follows that if A^{-1} exists,

$$x = (A^{-1}A)x = A^{-1}(Ax) = A^{-1}y$$

thus solving the system. Now in the case that A^{-1} exists, there is a formula for A^{-1} given above. Using this formula,

$$x_i = \sum_{j=1}^n a_{ij}^{-1} y_j = \sum_{j=1}^n \frac{1}{\det(A)} \operatorname{cof}(A)_{ji} y_j.$$

By the formula for the expansion of a determinant along a column,

$$x_i = \frac{1}{\det(A)} \det \begin{pmatrix} * & \cdots & y_1 & \cdots & * \\ \vdots & & \vdots & & \vdots \\ * & \cdots & y_n & \cdots & * \end{pmatrix},$$

where here the i^{th} column of A is replaced with the column vector, $(y_1, \dots, y_n)^T$, and the determinant of this modified matrix is taken and divided by $\det(A)$. This formula is known as Cramer's rule.

Definition 1.9.20 *A matrix M , is upper triangular if $M_{ij} = 0$ whenever $i > j$. Thus such a matrix equals zero below the main diagonal, the entries of the form M_{ii} as shown.*

$$\begin{pmatrix} * & * & \cdots & * \\ 0 & * & \ddots & \vdots \\ \vdots & \ddots & \ddots & * \\ 0 & \cdots & 0 & * \end{pmatrix}$$

A lower triangular matrix is defined similarly as a matrix for which all entries above the main diagonal are equal to zero.

With this definition, here is a simple corollary of Theorem 1.9.17.

Corollary 1.9.21 *Let M be an upper (lower) triangular matrix. Then $\det(M)$ is obtained by taking the product of the entries on the main diagonal.*

1.9.8 Rank of a Matrix

Definition 1.9.22 *A submatrix of a matrix A is the rectangular array of numbers obtained by deleting some rows and columns of A . Let A be an $m \times n$ matrix. The **determinant rank** of the matrix equals r where r is the largest number such that some $r \times r$ submatrix of A has a non zero determinant. The **row rank** is defined to be the dimension of the span of the rows. The **column rank** is defined to be the dimension of the span of the columns.*

Theorem 1.9.23 *If A , an $m \times n$ matrix has determinant rank r , then there exist r rows of the matrix such that every other row is a linear combination of these r rows.*

Proof: Suppose the determinant rank of $A = (a_{ij})$ equals r . Thus some $r \times r$ submatrix has non zero determinant and there is no larger square submatrix which has non zero determinant. Suppose such a submatrix is determined by the r columns whose indices are

$$j_1 < \cdots < j_r$$

and the r rows whose indices are

$$i_1 < \cdots < i_r$$

I want to show that every row is a linear combination of these rows. Consider the l^{th} row and let p be an index between 1 and n . Form the following $(r+1) \times (r+1)$ matrix

$$\begin{pmatrix} a_{i_1 j_1} & \cdots & a_{i_1 j_r} & a_{i_1 p} \\ \vdots & & \vdots & \vdots \\ a_{i_r j_1} & \cdots & a_{i_r j_r} & a_{i_r p} \\ a_{l j_1} & \cdots & a_{l j_r} & a_{l p} \end{pmatrix}$$

Of course you can assume $l \notin \{i_1, \dots, i_r\}$ because there is nothing to prove if the l^{th} row is one of the chosen ones. The above matrix has determinant 0. This is because if $p \notin \{j_1, \dots, j_r\}$ then the above would be a submatrix of A which is too large to have non zero determinant. On the other hand, if $p \in \{j_1, \dots, j_r\}$ then the above matrix has two columns which are equal so its determinant is still 0.

Expand the determinant of the above matrix along the last column. Let C_k denote the cofactor associated with the entry $a_{i_k p}$. This is not dependent on the choice of p . Remember, you delete the column and the row the entry is in and take the determinant of what is left and multiply by -1 raised to an appropriate power. Let C denote the cofactor associated with $a_{l p}$. This is given to be nonzero, it being the determinant of the matrix $r \times r$ matrix in the upper left corner. Thus $0 = a_{l p} C + \sum_{k=1}^r C_k a_{i_k p}$ which implies $a_{l p} = \sum_{k=1}^r \frac{-C_k}{C} a_{i_k p} \equiv \sum_{k=1}^r m_k a_{i_k p}$. Since this is true for every p and since m_k does not depend on p , this has shown the l^{th} row is a linear combination of the i_1, i_2, \dots, i_r rows. ■

Corollary 1.9.24 *The determinant rank equals the row rank.*

Proof: From Theorem 1.9.23, every row is in the span of r rows where r is the determinant rank. Therefore, the row rank (dimension of the span of the rows) is no larger than the determinant rank. Could the row rank be smaller than the determinant rank? If so, it follows from Theorem 1.9.23 that there exist p rows for $p < r \equiv$ determinant rank, such that the span of these p rows equals the row space. But then you could consider the $r \times r$ sub matrix which determines the determinant rank and it would follow that each of these rows would be in the span of the restrictions of the p rows just mentioned. By Theorem 4.2.3, the exchange theorem, the rows of this sub matrix would not be linearly independent and so some row is a linear combination of the others. By Corollary 1.9.11 the determinant would be 0, a contradiction. ■

Corollary 1.9.25 *If A has determinant rank r , then there exist r columns of the matrix such that every other column is a linear combination of these r columns. Also the column rank equals the determinant rank.*

Proof: This follows from the above by considering A^T . The rows of A^T are the columns of A and the determinant rank of A^T and A are the same. Therefore, from Corollary 1.9.24, column rank of $A =$ row rank of $A^T =$ determinant rank of $A^T =$ determinant rank of A . ■

The following theorem is of fundamental importance and ties together many of the ideas presented above.

Theorem 1.9.26 *Let A be an $n \times n$ matrix. Then the following are equivalent.*

1. $\det(A) = 0$.
2. A, A^T are not one to one.
3. A is not onto.

Proof: Suppose $\det(A) = 0$. Then the determinant rank of $A = r < n$. Therefore, there exist r columns such that every other column is a linear combination of these columns by Theorem 1.9.23. In particular, it follows that for some m , the m^{th} column is a linear combination of all the others. Thus letting $A = \begin{pmatrix} \mathbf{a}_1 & \cdots & \mathbf{a}_m & \cdots & \mathbf{a}_n \end{pmatrix}$ where the columns are denoted by \mathbf{a}_i , there exists scalars α_i such that $\mathbf{a}_m = \sum_{k \neq m} \alpha_k \mathbf{a}_k$. Now consider the column vector, $\mathbf{x} \equiv \begin{pmatrix} \alpha_1 & \cdots & -1 & \cdots & \alpha_n \end{pmatrix}^T$. Then $A\mathbf{x} = -\mathbf{a}_m + \sum_{k \neq m} \alpha_k \mathbf{a}_k = \mathbf{0}$. Since also $A\mathbf{0} = \mathbf{0}$, it follows A is not one to one. Similarly, A^T is not one to one by the same argument applied to A^T . This verifies that 1.) implies 2.).

Now suppose 2.). Then since A^T is not one to one, it follows there exists $\mathbf{x} \neq \mathbf{0}$ such that $A^T \mathbf{x} = \mathbf{0}$. Taking the transpose of both sides yields $\mathbf{x}^T A = \mathbf{0}^T$ where the $\mathbf{0}^T$ is a $1 \times n$ matrix or row vector. Now if $A\mathbf{y} = \mathbf{x}$, then $|\mathbf{x}|^2 = \mathbf{x}^T (A\mathbf{y}) = (\mathbf{x}^T A) \mathbf{y} = \mathbf{0} \mathbf{y} = 0$ contrary to $\mathbf{x} \neq \mathbf{0}$. Consequently there can be no \mathbf{y} such that $A\mathbf{y} = \mathbf{x}$ and so A is not onto. This shows that 2.) implies 3.).

Finally, suppose 3.). If 1.) does not hold, then $\det(A) \neq 0$ but then from Theorem 1.9.18 A^{-1} exists and so for every $\mathbf{y} \in \mathbb{F}^n$ there exists a unique $\mathbf{x} \in \mathbb{F}^n$ such that $A\mathbf{x} = \mathbf{y}$. In fact $\mathbf{x} = A^{-1}\mathbf{y}$. Thus A would be onto contrary to 3.). This shows 3.) implies 1.). ■

Corollary 1.9.27 *Let A be an $n \times n$ matrix. Then the following are equivalent.*

1. $\det(A) \neq 0$.

2. A and A^T are one to one.

3. A is onto.

Proof: This follows immediately from the above theorem.

1.9.9 An Identity of Cauchy

Theorem 1.9.28 *Both the left and the right sides in the following yield the same polynomial in the variables a_i, b_i for $i \leq n$.*

$$\prod_{i,j} (a_i + b_j) \begin{vmatrix} \frac{1}{a_1+b_1} & \cdots & \frac{1}{a_1+b_n} \\ \vdots & & \vdots \\ \frac{1}{a_n+b_1} & \cdots & \frac{1}{a_n+b_n} \end{vmatrix} = \prod_{j < i} (a_i - a_j) (b_i - b_j). \quad (1.26)$$

Proof: The theorem is true if $n = 2$. This follows from some computations. Suppose it is true for $n - 1$, $n \geq 3$.

$$\begin{aligned} & \begin{vmatrix} \frac{1}{a_1+b_1} & \frac{1}{a_1+b_2} & \cdots & \frac{1}{a_1+b_n} \\ \vdots & \vdots & \cdots & \vdots \\ \frac{1}{a_{n-1}+b_1} & \frac{1}{a_{n-1}+b_2} & \cdots & \frac{1}{a_{n-1}+b_n} \\ \frac{1}{a_n+b_1} & \frac{1}{a_n+b_2} & \cdots & \frac{1}{a_n+b_n} \end{vmatrix} \\ &= \begin{vmatrix} \frac{a_n-a_1}{(a_1+b_1)(b_1+a_n)} & \frac{a_n-a_1}{(a_1+b_2)(b_2+a_n)} & \cdots & \frac{a_n-a_1}{(a_1+b_n)(b_n+a_n)} \\ \vdots & \vdots & \cdots & \vdots \\ \frac{a_n-a_{n-1}}{(a_{n-1}+b_1)(b_1+a_n)} & \frac{a_n-a_{n-1}}{(b_2+a_n)(b_2+a_{n-1})} & \cdots & \frac{a_n-a_{n-1}}{(a_n+b_n)(b_n+a_{n-1})} \\ \frac{1}{a_n+b_1} & \frac{1}{a_n+b_2} & \cdots & \frac{1}{a_n+b_n} \end{vmatrix} \end{aligned}$$

Continuing to use the multilinear properties of determinants, this equals

$$\begin{vmatrix} \frac{1}{(a_1+b_1)(b_1+a_n)} & \frac{1}{(a_1+b_2)(b_2+a_n)} & \cdots & \frac{1}{(a_1+b_n)(b_n+a_n)} \\ \vdots & \vdots & \cdots & \vdots \\ \frac{1}{(a_{n-1}+b_1)(b_1+a_n)} & \frac{1}{(b_2+a_n)(b_2+a_{n-1})} & \cdots & \frac{1}{(a_n+b_n)(b_n+a_{n-1})} \\ \frac{1}{a_n+b_1} & \frac{1}{a_n+b_2} & \cdots & \frac{1}{a_n+b_n} \end{vmatrix} \prod_{k=1}^{n-1} (a_n - a_k)$$

and this equals

$$\begin{vmatrix} \frac{1}{(a_1+b_1)} & \frac{1}{(a_1+b_2)} & \cdots & \frac{1}{(a_1+b_n)} \\ \vdots & \vdots & \cdots & \vdots \\ \frac{1}{(a_{n-1}+b_1)} & \frac{1}{(b_2+a_{n-1})} & \cdots & \frac{1}{(b_n+a_{n-1})} \\ 1 & 1 & \cdots & 1 \end{vmatrix} \frac{\prod_{k=1}^{n-1} (a_n - a_k)}{\prod_{k=1}^n (a_n + b_k)}$$

Now take -1 times the last column and add to each previous column. Thus it equals

$$\begin{vmatrix} \frac{b_n-b_1}{(a_1+b_1)(a_1+b_n)} & \frac{b_n-b_2}{(a_1+b_2)(a_1+b_n)} & \cdots & \frac{1}{(a_1+b_n)} \\ \vdots & \vdots & \cdots & \vdots \\ \frac{b_n-b_1}{(b_1+a_{n-1})(b_n+a_{n-1})} & \frac{b_n-b_2}{(b_2+a_{n-1})(b_n+a_{n-1})} & \cdots & \frac{1}{(a_{n-1}+b_n)} \\ 0 & 0 & \cdots & 1 \end{vmatrix} \frac{\prod_{k=1}^{n-1} (a_n - a_k)}{\prod_{k=1}^n (a_n + b_k)}$$

Now continue simplifying using the multilinear property of the determinant.

$$\begin{vmatrix} \frac{1}{(a_1+b_1)} & \frac{1}{(a_1+b_2)} & \cdots & 1 \\ \vdots & \vdots & \cdots & \vdots \\ \frac{1}{(b_1+a_{n-1})} & \frac{1}{(b_2+a_{n-1})} & \cdots & 1 \\ 0 & 0 & \cdots & 1 \end{vmatrix} \frac{\prod_{k=1}^{n-1} (a_n - a_k) \prod_{k=1}^{n-1} (b_n - b_k)}{\prod_{k=1}^n (a_n + b_k) \prod_{k=1}^{n-1} (a_k + b_n)}$$

Expanding along the bottom row, what has just resulted is

$$\begin{vmatrix} \frac{1}{a_1+b_1} & \cdots & \frac{1}{a_1+b_{n-1}} \\ \vdots & \cdots & \vdots \\ \frac{1}{a_{n-1}+b_1} & \cdots & \frac{1}{a_{n-1}+b_{n-1}} \end{vmatrix} \frac{\prod_{k=1}^{n-1} (a_n - a_k) \prod_{k=1}^{n-1} (b_n - b_k)}{\prod_{k=1}^n (a_n + b_k) \prod_{k=1}^{n-1} (a_k + b_n)}$$

By induction this equals

$$\begin{aligned} & \frac{\prod_{j < i \leq n-1} (a_i - a_j) (b_i - b_j) \prod_{k=1}^{n-1} (a_n - a_k) \prod_{k=1}^{n-1} (b_n - b_k)}{\prod_{i,j \leq n-1} (a_i + b_j) \prod_{k=1}^n (a_n + b_k) \prod_{k=1}^{n-1} (a_k + b_n)} \\ &= \frac{\prod_{j < i \leq n} (a_i - a_j) (b_i - b_j)}{\prod_{i,j \leq n} (a_i + b_j)} \blacksquare \end{aligned}$$

1.10 The Cayley Hamilton Theorem

Definition 1.10.1 Let A be an $n \times n$ matrix. The characteristic polynomial is defined as

$$q_A(t) \equiv \det(tI - A)$$

and the solutions to $q_A(t) = 0$ are called eigenvalues. For A a matrix and $p(t) = t^n + a_{n-1}t^{n-1} + \cdots + a_1t + a_0$, denote by $p(A)$ the matrix defined by

$$p(A) \equiv A^n + a_{n-1}A^{n-1} + \cdots + a_1A + a_0I.$$

The explanation for the last term is that A^0 is interpreted as I , the identity matrix.

The Cayley Hamilton theorem states that every matrix satisfies its characteristic equation, that equation defined by $q_A(t) = 0$. It is one of the most important theorems in linear algebra². The proof in this section is not the most general proof, but works well when the field of scalars is \mathbb{R} or \mathbb{C} . The following lemma will help with its proof.

Lemma 1.10.2 Suppose for all $|\lambda|$ large enough,

$$A_0 + A_1\lambda + \cdots + A_m\lambda^m = 0,$$

where the A_i are $n \times n$ matrices. Then each $A_i = 0$.

²A special case was first proved by Hamilton in 1853. The general case was announced by Cayley some time later and a proof was given by Frobenius in 1878.

Proof: Suppose some $A_i \neq 0$. Let p be the largest index of those which are non zero. Then multiply by λ^{-p} .

$$A_0\lambda^{-p} + A_1\lambda^{-p+1} + \cdots + A_{p-1}\lambda^{-1} + A_p = 0$$

Now let $\lambda \rightarrow \infty$. Thus $A_p = 0$ after all. Hence each $A_i = 0$. ■

With the lemma, here is a simple corollary.

Corollary 1.10.3 *Let A_i and B_i be $n \times n$ matrices and suppose*

$$A_0 + A_1\lambda + \cdots + A_m\lambda^m = B_0 + B_1\lambda + \cdots + B_m\lambda^m$$

for all $|\lambda|$ large enough. Then $A_i = B_i$ for all i . If $A_i = B_i$ for each A_i, B_i then one can substitute an $n \times n$ matrix M for λ and the identity will continue to hold.

Proof: Subtract and use the result of the lemma. The last claim is obvious by matching terms. ■

With this preparation, here is a relatively easy proof of the Cayley Hamilton theorem.

Theorem 1.10.4 *Let A be an $n \times n$ matrix and let $q(\lambda) \equiv \det(\lambda I - A)$ be the characteristic polynomial. Then $q(A) = 0$.*

Proof: Let $C(\lambda)$ equal the transpose of the cofactor matrix of $(\lambda I - A)$ for $|\lambda|$ large. (If $|\lambda|$ is large enough, then λ cannot be in the finite list of eigenvalues of A and so for such λ , $(\lambda I - A)^{-1}$ exists.) Therefore, by Theorem 1.9.18

$$C(\lambda) = q(\lambda)(\lambda I - A)^{-1}.$$

Say

$$q(\lambda) = a_0 + a_1\lambda + \cdots + \lambda^n$$

Note that each entry in $C(\lambda)$ is a polynomial in λ having degree no more than $n-1$. For example, you might have something like

$$\begin{aligned} C(\lambda) &= \begin{pmatrix} \lambda^2 - 6\lambda + 9 & 3 - \lambda & 0 \\ 2\lambda - 6 & \lambda^2 - 3\lambda & 0 \\ \lambda - 1 & \lambda - 1 & \lambda^2 - 3\lambda + 2 \end{pmatrix} \\ &= \begin{pmatrix} 9 & 3 & 0 \\ -6 & 0 & 0 \\ -1 & -1 & 2 \end{pmatrix} + \lambda \begin{pmatrix} -6 & -1 & 0 \\ 2 & -3 & 0 \\ 1 & 1 & -3 \end{pmatrix} + \lambda^2 \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \end{aligned}$$

Therefore, collecting the terms in the general case,

$$C(\lambda) = C_0 + C_1\lambda + \cdots + C_{n-1}\lambda^{n-1}$$

for C_j some $n \times n$ matrix. Then

$$C(\lambda)(\lambda I - A) = (C_0 + C_1\lambda + \cdots + C_{n-1}\lambda^{n-1})(\lambda I - A) = q(\lambda)I$$

Then multiplying out the middle term, it follows that for all $|\lambda|$ sufficiently large,

$$\begin{aligned} a_0I + a_1I\lambda + \cdots + I\lambda^n &= C_0\lambda + C_1\lambda^2 + \cdots + C_{n-1}\lambda^n \\ &\quad - [C_0A + C_1A\lambda + \cdots + C_{n-1}A\lambda^{n-1}] \\ &= -C_0A + (C_0 - C_1A)\lambda + (C_1 - C_2A)\lambda^2 + \cdots + (C_{n-2} - C_{n-1}A)\lambda^{n-1} + C_{n-1}\lambda^n \end{aligned}$$

Then, using Corollary 1.10.3, one can replace λ on both sides with A . Then the right side is seen to equal 0. Hence the left side, $q(A)I$ is also equal to 0. ■

Part I

Topology, Continuity, Algebra, Derivatives

Chapter 2

Some Basic Topics

This chapter contains basic definitions and a few fundamental theorems which will be used throughout the book whenever convenient.

2.1 Basic Definitions

A set is a collection of things called elements of the set. For example, the set of integers, the collection of signed whole numbers such as $1, 2, -4$, etc. This set whose existence will be assumed is denoted by \mathbb{Z} . Other sets could be the set of people in a family or the set of donuts in a display case at the store. Sometimes parentheses, $\{ \}$ specify a set by listing the things which are in the set between the parentheses. For example the set of integers between -1 and 2 , including these numbers could be denoted as $\{-1, 0, 1, 2\}$. The notation signifying x is an element of a set S , is written as $x \in S$. Thus, $1 \in \{-1, 0, 1, 2, 3\}$. Here are some axioms about sets. Axioms are statements which are accepted, not proved.

Axiom 2.1.1 *Two sets are equal if and only if they have the same elements.*

Axiom 2.1.2 *To every set, A , and to every condition $S(x)$ there corresponds a set, B , whose elements are exactly those elements x of A for which $S(x)$ holds.*

Axiom 2.1.3 *For every collection of sets there exists a set that contains all the elements that belong to at least one set of the given collection. (You can take the union of a bunch of sets.)*

Axiom 2.1.4 *The Cartesian product of a nonempty family of nonempty sets is nonempty.*

Axiom 2.1.5 *If A is a set there exists a set, $\mathcal{P}(A)$ such that $\mathcal{P}(A)$ is the set of all subsets of A . This is called the power set.*

These axioms are referred to as the axiom of extension, axiom of specification, axiom of unions, axiom of choice, and axiom of powers respectively.

It seems fairly clear you should want to believe in the axiom of extension. It is merely saying, for example, that $\{1, 2, 3\} = \{2, 3, 1\}$ since these two sets have the same elements in them. Similarly, it would seem you should be able to specify a new set from a given set using some “condition” which can be used as a test to determine whether the element in question is in the set. For example, the set of all integers which are multiples of 2. This set could be specified as follows.

$$\{x \in \mathbb{Z} : x = 2y \text{ for some } y \in \mathbb{Z}\}.$$

In this notation, the colon is read as “such that” and in this case the condition is being a multiple of 2.

Another example of political interest, could be the set of all judges who are not judicial activists. I think you can see this last is not a very precise condition since there is no way to determine to everyone’s satisfaction whether a given judge is an activist. Also, **just because something is grammatically correct does not mean it makes any sense**. For example consider the following nonsense.

$$S = \{x \in \text{set of dogs} : \text{it is colder in the mountains than in the winter}\}.$$

So what is a condition?

We will leave these sorts of considerations and assume our conditions make sense, whatever that means. The axiom of unions states that for any collection of sets, there is a set consisting of all the elements in each of the sets in the collection. Of course this is also open to further consideration. What is a collection? Maybe it would be better to say “set of sets” or, given a set whose elements are sets there exists a set whose elements consist of exactly those things which are elements of at least one of these sets. If \mathcal{S} is such a set whose elements are sets,

$$\cup\{A : A \in \mathcal{S}\} \text{ or } \cup\mathcal{S}$$

signify this union.

Something is in the Cartesian product of a set or “family” of sets if it consists of a single thing taken from each set in the family. Thus $(1, 2, 3) \in \{1, 4, 2\} \times \{1, 2, 7\} \times \{4, 3, 7, 9\}$ because it consists of exactly one element from each of the sets which are separated by \times . Also, this is the notation for the Cartesian product of finitely many sets. If \mathcal{S} is a set whose elements are sets, $\prod_{A \in \mathcal{S}} A$ signifies the Cartesian product.

The Cartesian product is the set of choice functions, a choice function being a function which selects exactly one element of each set of \mathcal{S} . You may think the axiom of choice, stating that the Cartesian product of a nonempty family of nonempty sets is nonempty, is innocuous but there was a time when many mathematicians were ready to throw it out because it implies things which are very hard to believe, things which never happen without the axiom of choice.

A is a subset of B , written $A \subseteq B$, if every element of A is also an element of B . This can also be written as $B \supseteq A$. A is a proper subset of B , written $A \subset B$ or $B \supset A$ if A is a subset of B but A is not equal to B , $A \neq B$. $A \cap B$ denotes the intersection of the two sets, A and B and it means the set of elements of A which are also elements of B . The axiom of specification shows this is a set. The empty set is the set which has no elements in it, denoted as \emptyset . $A \cup B$ denotes the union of the two sets, A and B and it means the set of all elements which are in either of the sets. It is a set because of the axiom of unions.

The complement of a set, (the set of things which are not in the given set) must be taken with respect to a given set called the universal set which is a set which contains the one whose complement is being taken. Thus, the complement of A , denoted as A^C (or more precisely as $X \setminus A$) is a set obtained from using the axiom of specification to write

$$A^C \equiv \{x \in X : x \notin A\}$$

The symbol \notin means: “is not an element of”. Note the axiom of specification takes place relative to a given set. Without this universal set it makes no sense to use the axiom of specification to obtain the complement.

Words such as “all” or “there exists” are called quantifiers and they must be understood relative to some given set. For example, the set of all integers larger than 3. Or there exists an integer larger than 7. Such statements have to do with a given set, in this case the integers. Failure to have a reference set when quantifiers are used turns out to be illogical even though such usage may be grammatically correct. Quantifiers are used often enough that there are symbols for them. The symbol \forall is read as “for all” or “for every” and the symbol \exists is read as “there exists”. Thus $\forall \exists E$ could mean for every upside down A there exists a backwards E .

DeMorgan’s laws are very useful in mathematics. Let \mathcal{S} be a set of sets each of which

is contained in some universal set, U . Then

$$\cup \{A^C : A \in \mathcal{S}\} = (\cap \{A : A \in \mathcal{S}\})^C$$

and

$$\cap \{A^C : A \in \mathcal{S}\} = (\cup \{A : A \in \mathcal{S}\})^C.$$

These laws follow directly from the definitions. Also following directly from the definitions are:

Let \mathcal{S} be a set of sets then

$$B \cup \cup \{A : A \in \mathcal{S}\} = \cup \{B \cup A : A \in \mathcal{S}\}.$$

and: Let \mathcal{S} be a set of sets show

$$B \cap \cup \{A : A \in \mathcal{S}\} = \cup \{B \cap A : A \in \mathcal{S}\}.$$

Unfortunately, there is no single universal set which can be used for all sets. Here is why: Suppose there were. Call it S . Then you could consider A the set of all elements of S which are not elements of themselves, this from the axiom of specification. If A is an element of itself, then it fails to qualify for inclusion in A . Therefore, it must not be an element of itself. However, if this is so, it qualifies for inclusion in A so it is an element of itself and so this can't be true either. Thus the most basic of conditions you could imagine, that of being an element of, is meaningless and so allowing such a set causes the whole theory to be meaningless. The solution is to not allow a universal set. As mentioned by Halmos in Naive set theory, "Nothing contains everything". Always beware of statements involving quantifiers wherever they occur, even this one. This little observation described above is due to Bertrand Russell and is called Russell's paradox.

2.2 The Schroder Bernstein Theorem

It is very important to be able to compare the size of sets in a rational way. The most useful theorem in this context is the Schroder Bernstein theorem which is the main result to be presented in this section. The Cartesian product is discussed above. The next definition reviews this and defines the concept of a function.

Definition 2.2.1 *Let X and Y be sets.*

$$X \times Y \equiv \{(x, y) : x \in X \text{ and } y \in Y\}$$

A relation is defined to be a subset of $X \times Y$. A function f , also called a mapping, is a relation which has the property that if (x, y) and (x, y_1) are both elements of the f , then $y = y_1$. The domain of f is defined as

$$D(f) \equiv \{x : (x, y) \in f\},$$

written as $f : D(f) \rightarrow Y$. Another notation which is used is the following

$$f^{-1}(y) \equiv \{x \in D(f) : f(x) = y\}$$

This is called the inverse image.

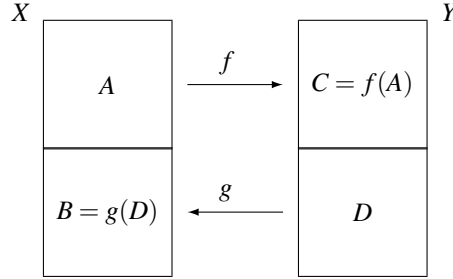
It is probably safe to say that most people do not think of functions as a type of relation which is a subset of the Cartesian product of two sets. A function is like a machine which takes inputs, x and makes them into a unique output, $f(x)$. Of course, that is what the above definition says with more precision. An ordered pair, (x, y) which is an element of the function or mapping has an input, x and a unique output y , denoted as $f(x)$ while the name of the function is f . “mapping” is often a noun meaning function. However, it also is a verb as in “ f is mapping A to B ”. That which a function is thought of as doing is also referred to using the word “maps” as in: f maps X to Y . However, a set of functions may be called a set of maps so this word might also be used as the plural of a noun. There is no help for it. You just have to suffer with this nonsense.

The following theorem which is interesting for its own sake will be used to prove the Schroder Bernstein theorem.

Theorem 2.2.2 *Let $f : X \rightarrow Y$ and $g : Y \rightarrow X$ be two functions. Then there exist sets A, B, C, D , such that*

$$A \cup B = X, C \cup D = Y, A \cap B = \emptyset, C \cap D = \emptyset, \\ f(A) = C, g(D) = B.$$

The following picture illustrates the conclusion of this theorem.



Proof: Consider the empty set, $\emptyset \subseteq X$. If $y \in Y \setminus f(\emptyset)$, then $g(y) \notin \emptyset$ because \emptyset has no elements. Also, if A, B, C , and D are as described above, A also would have this same property that the empty set has. However, A is probably larger. Therefore, say $A_0 \subseteq X$ satisfies \mathcal{P} if whenever $y \in Y \setminus f(A_0)$, $g(y) \notin A_0$.

$$\mathcal{A} \equiv \{A_0 \subseteq X : A_0 \text{ satisfies } \mathcal{P}\}.$$

Let $A = \cup \mathcal{A}$. If $y \in Y \setminus f(A)$, then for each $A_0 \in \mathcal{A}$, $y \in Y \setminus f(A_0)$ and so $g(y) \notin A_0$. Since $g(y) \notin A_0$ for all $A_0 \in \mathcal{A}$, it follows $g(y) \notin A$. Hence A satisfies \mathcal{P} and is the largest subset of X which does so. Now define

$$C \equiv f(A), D \equiv Y \setminus C, B \equiv X \setminus A.$$

It only remains to verify that $g(D) = B$. It was just shown that $g(D) \subseteq B$.

Suppose $x \in B = X \setminus A$. Then $A \cup \{x\}$ does not satisfy \mathcal{P} and so there exists $y \in Y \setminus f(A \cup \{x\}) \subseteq D$ such that $g(y) \in A \cup \{x\}$. But $y \notin f(A)$ and so since A satisfies \mathcal{P} , it follows $g(y) \notin A$. Hence $g(y) = x$ and so $x \in g(D)$. Hence $g(D) = B$. ■

Theorem 2.2.3 (Schroder Bernstein) *If $f : X \rightarrow Y$ and $g : Y \rightarrow X$ are one to one, then there exists $h : X \rightarrow Y$ which is one to one and onto.*

Proof: Let A, B, C, D be the sets of Theorem 2.2.2 and define

$$h(x) \equiv \begin{cases} f(x) & \text{if } x \in A \\ g^{-1}(x) & \text{if } x \in B \end{cases}$$

Then h is the desired one to one and onto mapping. ■

Recall that the Cartesian product may be considered as the collection of choice functions.

Definition 2.2.4 Let I be a set and let X_i be a set for each $i \in I$. f is a choice function written as $f \in \prod_{i \in I} X_i$ if $f(i) \in X_i$ for each $i \in I$.

The axiom of choice says that if $X_i \neq \emptyset$ for each $i \in I$, for I a set, then

$$\prod_{i \in I} X_i \neq \emptyset.$$

Sometimes the two functions, f and g are onto but not one to one. It turns out that with the axiom of choice, a similar conclusion to the above may be obtained.

Corollary 2.2.5 If $f : X \rightarrow Y$ is onto and $g : Y \rightarrow X$ is onto, then there exists $h : X \rightarrow Y$ which is one to one and onto.

Proof: For each $y \in Y$, $f^{-1}(y) \equiv \{x \in X : f(x) = y\} \neq \emptyset$. Therefore, by the axiom of choice, there exists $f_0^{-1} \in \prod_{y \in Y} f^{-1}(y)$ which is the same as saying that for each $y \in Y$, $f_0^{-1}(y) \in f^{-1}(y)$. Similarly, there exists $g_0^{-1}(x) \in g^{-1}(x)$ for all $x \in X$. Then f_0^{-1} is one to one because if $f_0^{-1}(y_1) = f_0^{-1}(y_2)$, then

$$y_1 = f(f_0^{-1}(y_1)) = f(f_0^{-1}(y_2)) = y_2.$$

Similarly g_0^{-1} is one to one. Therefore, by the Schroder Bernstein theorem, there exists $h : X \rightarrow Y$ which is one to one and onto. ■

Definition 2.2.6 A set S , is finite if there exists a natural number n and a map θ which maps $\{1, \dots, n\}$ one to one and onto S . S is infinite if it is not finite. A set S , is called countable if there exists a map θ mapping \mathbb{N} one to one and onto S . (When θ maps a set A to a set B , this will be written as $\theta : A \rightarrow B$ in the future.) Here $\mathbb{N} \equiv \{1, 2, \dots\}$, the natural numbers. S is at most countable if there exists a map $\theta : \mathbb{N} \rightarrow S$ which is onto.

The property of being at most countable is often referred to as being countable because the question of interest is normally whether one can list all elements of the set, designating a first, second, third etc. in such a way as to give each element of the set a natural number. The possibility that a single element of the set may be counted more than once is often not important.

Theorem 2.2.7 If X and Y are both at most countable, then $X \times Y$ is also at most countable. If either X or Y is countable, then $X \times Y$ is also countable.

Proof: It is given that there exists a mapping $\eta : \mathbb{N} \rightarrow X$ which is onto. Define $\eta(i) \equiv x_i$ and consider X as the set $\{x_1, x_2, x_3, \dots\}$. Similarly, consider Y as the set $\{y_1, y_2, y_3, \dots\}$. It follows the elements of $X \times Y$ are included in the following rectangular array.

$$\begin{array}{ccccccc} (x_1, y_1) & (x_1, y_2) & (x_1, y_3) & \cdots & \leftarrow & \text{Those which have } x_1 \text{ in first slot.} \\ (x_2, y_1) & (x_2, y_2) & (x_2, y_3) & \cdots & \leftarrow & \text{Those which have } x_2 \text{ in first slot.} \\ (x_3, y_1) & (x_3, y_2) & (x_3, y_3) & \cdots & \leftarrow & \text{Those which have } x_3 \text{ in first slot.} \\ \vdots & \vdots & \vdots & & & \vdots \end{array}$$

Follow a path through this array as follows.

$$\begin{array}{ccccc} (x_1, y_1) & \rightarrow & (x_1, y_2) & & (x_1, y_3) \rightarrow \\ & \swarrow & & \nearrow & \\ (x_2, y_1) & & (x_2, y_2) & & \\ \downarrow & \nearrow & & & \\ (x_3, y_1) & & & & \end{array}$$

Thus the first element of $X \times Y$ is (x_1, y_1) , the second element of $X \times Y$ is (x_1, y_2) , the third element of $X \times Y$ is (x_2, y_1) etc. This assigns a number from \mathbb{N} to each element of $X \times Y$. Thus $X \times Y$ is at most countable.

It remains to show the last claim. Suppose without loss of generality that X is countable. Then there exists $\alpha : \mathbb{N} \rightarrow X$ which is one to one and onto. Let $\beta : X \times Y \rightarrow \mathbb{N}$ be defined by $\beta((x, y)) \equiv \alpha^{-1}(x)$. Thus β is onto \mathbb{N} . By the first part there exists a function from \mathbb{N} onto $X \times Y$. Therefore, by Corollary 2.2.5, there exists a one to one and onto mapping from $X \times Y$ to \mathbb{N} . ■

Theorem 2.2.8 *If X and Y are at most countable, then $X \cup Y$ is at most countable. If either X or Y are countable, then $X \cup Y$ is countable.*

Proof: As in the preceding theorem,

$$X = \{x_1, x_2, x_3, \dots\}$$

and

$$Y = \{y_1, y_2, y_3, \dots\}.$$

Consider the following array consisting of $X \cup Y$ and path through it.

$$\begin{array}{ccccc} x_1 & \rightarrow & x_2 & & x_3 \rightarrow \\ & \swarrow & & \nearrow & \\ y_1 & \rightarrow & y_2 & & \end{array}$$

Thus the first element of $X \cup Y$ is x_1 , the second is x_2 the third is y_1 the fourth is y_2 etc.

Consider the second claim. By the first part, there is a map from \mathbb{N} onto $X \times Y$. Suppose without loss of generality that X is countable and $\alpha : \mathbb{N} \rightarrow X$ is one to one and onto. Then define $\beta(y) \equiv 1$, for all $y \in Y$, and $\beta(x) \equiv \alpha^{-1}(x)$. Thus, β maps $X \times Y$ onto \mathbb{N} and this shows there exist two onto maps, one mapping $X \cup Y$ onto \mathbb{N} and the other mapping \mathbb{N} onto $X \cup Y$. Then Corollary 2.2.5 yields the conclusion. ■

Note that by induction this shows that if you have any finite set whose elements are countable sets, then the union of these is countable.

2.3 Equivalence Relations

There are many ways to compare elements of a set other than to say two elements are equal or the same. For example, in the set of people let two people be equivalent if they have the same weight. This would not be saying they were the same person, just that they weighed the same. Often such relations involve considering one characteristic of the elements of a set and then saying the two elements are equivalent if they are the same as far as the given characteristic is concerned.

Definition 2.3.1 *Let S be a set. \sim is an equivalence relation on S if it satisfies the following axioms.*

1. $x \sim x$ for all $x \in S$. (Reflexive)
2. If $x \sim y$ then $y \sim x$. (Symmetric)
3. If $x \sim y$ and $y \sim z$, then $x \sim z$. (Transitive)

Definition 2.3.2 $[x]$ denotes the set of all elements of S which are equivalent to x and $[x]$ is called the equivalence class determined by x or just the equivalence class of x .

With the above definition one can prove the following simple theorem.

Theorem 2.3.3 *Let \sim be an equivalence relation defined on a set, S and let \mathcal{H} denote the set of equivalence classes. Then if $[x]$ and $[y]$ are two of these equivalence classes, either $x \sim y$ and $[x] = [y]$ or it is not true that $x \sim y$ and $[x] \cap [y] = \emptyset$.*

2.4 sup and inf

It is assumed in all that is done that \mathbb{R} is complete. There are two ways to describe completeness of \mathbb{R} . One is to say that every bounded set has a least upper bound and a greatest lower bound. The other is to say that every Cauchy sequence converges. These two equivalent notions of completeness will be taken as given. Cauchy sequences are discussed a little later.

The symbol, \mathbb{F} will mean either \mathbb{R} or \mathbb{C} . The symbol $[-\infty, \infty]$ will mean all real numbers along with $+\infty$ and $-\infty$ which are points which we pretend are at the right and left ends of the real line respectively. The inclusion of these make believe points makes the statement of certain theorems less trouble.

Definition 2.4.1 *For $A \subseteq [-\infty, \infty]$, $A \neq \emptyset$ $\sup A$ is defined as the least upper bound in case A is bounded above by a real number and equals ∞ if A is not bounded above. Similarly $\inf A$ is defined to equal the greatest lower bound in case A is bounded below by a real number and equals $-\infty$ in case A is not bounded below.*

Lemma 2.4.2 *If $\{A_n\}$ is an increasing sequence in $[-\infty, \infty]$, then*

$$\sup \{A_n : n \in \mathbb{N}\} = \lim_{n \rightarrow \infty} A_n.$$

Similarly, if $\{A_n\}$ is decreasing, then

$$\inf \{A_n : n \in \mathbb{N}\} = \lim_{n \rightarrow \infty} A_n.$$

Proof: Let $\sup(\{A_n : n \in \mathbb{N}\}) = r$. In the first case, suppose $r < \infty$. Then letting $\varepsilon > 0$ be given, there exists n such that $A_n \in (r - \varepsilon, r]$. Since $\{A_n\}$ is increasing, it follows if $m > n$, then $r - \varepsilon < A_n \leq A_m \leq r$ and so $\lim_{n \rightarrow \infty} A_n = r$ as claimed. In the case where $r = \infty$, then if a is a real number, there exists n such that $A_n > a$. Since $\{A_k\}$ is increasing, it follows that if $m > n$, $A_m > a$. But this is what is meant by $\lim_{n \rightarrow \infty} A_n = \infty$. The other case is that $r = -\infty$. But in this case, $A_n = -\infty$ for all n and so $\lim_{n \rightarrow \infty} A_n = -\infty$. The case where A_n is decreasing is entirely similar. ■

2.5 Double Series

Double series are of the form $\sum_{k=m}^{\infty} \sum_{j=m}^{\infty} a_{jk} \equiv \sum_{k=m}^{\infty} (\sum_{j=m}^{\infty} a_{jk})$. In other words, first sum on j yielding something which depends on k and then sum these. The major consideration for these double series is the question of when $\sum_{k=m}^{\infty} \sum_{j=m}^{\infty} a_{jk} = \sum_{j=m}^{\infty} \sum_{k=m}^{\infty} a_{jk}$. In other words, when does it make no difference which subscript is summed over first? In the case of finite sums there is no issue here. You can always write $\sum_{k=m}^M \sum_{j=m}^N a_{jk} = \sum_{j=m}^N \sum_{k=m}^M a_{jk}$ because addition is commutative. However, there are limits involved with infinite sums and the interchange in order of summation involves taking limits in a different order. Therefore, it is not always true that it is permissible to interchange the two sums. A general rule of thumb is this: If something involves changing the order in which two limits are taken, you may not do it without agonizing over the question. In general, limits foul up algebra and also introduce things which are counter intuitive. Here is an example. This example is a little technical. It is placed here just to prove conclusively there is a question which needs to be considered.

Example 2.5.1 Consider the following picture which depicts some of the ordered pairs (m, n) where m, n are positive integers.

$$\begin{array}{ccccc} & & & \vdots & \\ & 0 & 0 & c & 0 & -c \\ & 0 & c & 0 & -c & 0 \\ & b & 0 & -c & 0 & 0 & \cdots \\ & 0 & a & 0 & 0 & 0 \end{array}$$

The a, b, c are the values of a_{mn} . Thus $a_{nn} = 0$ for all $n \geq 1$, $a_{21} = a$, $a_{12} = b$, $a_{m(m+1)} = -c$ whenever $m > 1$, and $a_{m(m-1)} = c$ whenever $m > 2$. The numbers next to the point are the values of a_{mn} . You see $a_{nn} = 0$ for all n , $a_{21} = a$, $a_{12} = b$, $a_{mn} = c$ for (m, n) on the line $y = 1 + x$ whenever $m > 1$, and $a_{mn} = -c$ for all (m, n) on the line $y = x - 1$ whenever $m > 2$.

Then $\sum_{m=1}^{\infty} a_{mn} = a$ if $n = 1$, $\sum_{m=1}^{\infty} a_{mn} = b - c$ if $n = 2$ and if $n > 2$, $\sum_{m=1}^{\infty} a_{mn} = 0$. Therefore,

$$\sum_{n=1}^{\infty} \sum_{m=1}^{\infty} a_{mn} = a + b - c.$$

Next observe that $\sum_{n=1}^{\infty} a_{mn} = b$ if $m = 1$, $\sum_{n=1}^{\infty} a_{mn} = a + c$ if $m = 2$, and $\sum_{n=1}^{\infty} a_{mn} = 0$ if $m > 2$. Therefore,

$$\sum_{m=1}^{\infty} \sum_{n=1}^{\infty} a_{mn} = b + a + c$$

and so the two sums are different. Moreover, you can see that by assigning different values of a, b , and c , you can get an example for any two different numbers desired.

It turns out that if $a_{ij} \geq 0$ for all i, j , then you can always interchange the order of summation. This is shown next and is based on the following lemma. First, some notation should be discussed.

Definition 2.5.2 Let $f(a, b) \in [-\infty, \infty]$ for $a \in A$ and $b \in B$ where A, B are sets which means that $f(a, b)$ is either a number, ∞ , or $-\infty$. The symbol, $+\infty$ is interpreted as a point out at the end of the number line which is larger than every real number. Of course there is no such number. That is why it is called ∞ . The symbol, $-\infty$ is interpreted similarly. Then $\sup_{a \in A} f(a, b)$ means $\sup(S_b)$ where $S_b \equiv \{f(a, b) : a \in A\}$.

Unlike limits, you can take the sup in different orders.

Lemma 2.5.3 Let $f(a, b) \in [-\infty, \infty]$ for $a \in A$ and $b \in B$ where A, B are sets. Then

$$\sup_{a \in A} \sup_{b \in B} f(a, b) = \sup_{b \in B} \sup_{a \in A} f(a, b).$$

Proof: Note that for all a, b , $f(a, b) \leq \sup_{b \in B} \sup_{a \in A} f(a, b)$ and therefore, for all a , $\sup_{b \in B} f(a, b) \leq \sup_{b \in B} \sup_{a \in A} f(a, b)$. Therefore,

$$\sup_{a \in A} \sup_{b \in B} f(a, b) \leq \sup_{b \in B} \sup_{a \in A} f(a, b).$$

Repeat the same argument interchanging a and b , to get the conclusion of the lemma. ■

Theorem 2.5.4 Let $a_{ij} \geq 0$. Then $\sum_{i=1}^{\infty} \sum_{j=1}^{\infty} a_{ij} = \sum_{j=1}^{\infty} \sum_{i=1}^{\infty} a_{ij}$.

Proof: First note there is no trouble in defining these sums because the a_{ij} are all nonnegative. If a sum diverges, it only diverges to ∞ and so ∞ is the value of the sum. Next note that

$$\sum_{j=r}^{\infty} \sum_{i=r}^{\infty} a_{ij} \geq \sup_n \sum_{j=r}^n \sum_{i=r}^n a_{ij}$$

because for all j , $\sum_{i=r}^{\infty} a_{ij} \geq \sum_{i=r}^n a_{ij}$. Therefore,

$$\begin{aligned} \sum_{j=r}^{\infty} \sum_{i=r}^{\infty} a_{ij} &\geq \sup_n \sum_{j=r}^n \sum_{i=r}^n a_{ij} = \sup_n \lim_{m \rightarrow \infty} \sum_{j=r}^m \sum_{i=r}^n a_{ij} \\ &= \sup_n \lim_{m \rightarrow \infty} \sum_{i=r}^n \sum_{j=r}^m a_{ij} = \sup_n \sum_{i=r}^n \lim_{m \rightarrow \infty} \sum_{j=r}^m a_{ij} \\ &= \sup_n \sum_{i=r}^n \sum_{j=r}^{\infty} a_{ij} = \lim_{n \rightarrow \infty} \sum_{i=r}^n \sum_{j=r}^{\infty} a_{ij} = \sum_{i=r}^{\infty} \sum_{j=r}^{\infty} a_{ij} \end{aligned}$$

Interchanging the i and j in the above argument proves the theorem. ■

2.6 lim sup and lim inf

Sometimes the limit of a sequence does not exist. For example, if $a_n = (-1)^n$, then $\lim_{n \rightarrow \infty} a_n$ does not exist. This is because the terms of the sequence are a distance of 1 apart. Therefore there can't exist a single number such that all the terms of the sequence are ultimately within $1/4$ of that number. The nice thing about limsup and liminf is that they always exist. First here is a simple lemma and definition. First review the definition of inf and sup on Page 63 along with the simple properties of these things.

Definition 2.6.1 Denote by $[-\infty, \infty]$ the real line along with symbols ∞ and $-\infty$. It is understood that ∞ is larger than every real number and $-\infty$ is smaller than every real number. Then if $\{A_n\}$ is an increasing sequence of points of $[-\infty, \infty]$, $\lim_{n \rightarrow \infty} A_n$ equals ∞ if the only upper bound of the set $\{A_n\}$ is ∞ . If $\{A_n\}$ is bounded above by a real number, then $\lim_{n \rightarrow \infty} A_n$ is defined in the usual way and equals the least upper bound of $\{A_n\}$. If $\{A_n\}$ is a decreasing sequence of points of $[-\infty, \infty]$, $\lim_{n \rightarrow \infty} A_n$ equals $-\infty$ if the only lower bound of the sequence $\{A_n\}$ is $-\infty$. If $\{A_n\}$ is bounded below by a real number, then $\lim_{n \rightarrow \infty} A_n$ is defined in the usual way and equals the greatest lower bound of $\{A_n\}$. More simply, if $\{A_n\}$ is increasing, $\lim_{n \rightarrow \infty} A_n \equiv \sup \{A_n\}$ and if $\{A_n\}$ is decreasing then $\lim_{n \rightarrow \infty} A_n \equiv \inf \{A_n\}$.

Lemma 2.6.2 Let $\{a_n\}$ be a sequence of real numbers and let $U_n \equiv \sup \{a_k : k \geq n\}$. Then $\{U_n\}$ is a decreasing sequence. Also if $L_n \equiv \inf \{a_k : k \geq n\}$, then $\{L_n\}$ is an increasing sequence. Therefore, $\lim_{n \rightarrow \infty} L_n$ and $\lim_{n \rightarrow \infty} U_n$ both exist.

Proof: Let W_n be an upper bound for $\{a_k : k \geq n\}$. Then since these sets are getting smaller, it follows that for $m < n$, W_m is an upper bound for $\{a_k : k \geq n\}$. In particular if $W_m = U_m$, then U_m is an upper bound for $\{a_k : k \geq n\}$ and so U_m is at least as large as U_n , the least upper bound for $\{a_k : k \geq n\}$. The claim that $\{L_n\}$ is decreasing is similar. ■

From the lemma, the following definition makes sense.

Definition 2.6.3 Let $\{a_n\}$ be any sequence of points of $[-\infty, \infty]$

$$\limsup_{n \rightarrow \infty} a_n \equiv \lim_{n \rightarrow \infty} \sup \{a_k : k \geq n\}$$

$$\liminf_{n \rightarrow \infty} a_n \equiv \lim_{n \rightarrow \infty} \inf \{a_k : k \geq n\}.$$

Theorem 2.6.4 Suppose $\{a_n\}$ is a sequence of real numbers and also that both $\limsup_{n \rightarrow \infty} a_n, \liminf_{n \rightarrow \infty} a_n$ are real numbers. Then $\lim_{n \rightarrow \infty} a_n$ exists if and only if the two numbers are equal and in this case, the limit and the each of $\limsup_{n \rightarrow \infty} a_n, \liminf_{n \rightarrow \infty} a_n$ are equal.

Proof: First note that $\sup \{a_k : k \geq n\} \geq \inf \{a_k : k \geq n\}$ and so,

$$\limsup_{n \rightarrow \infty} a_n \equiv \lim_{n \rightarrow \infty} \sup \{a_k : k \geq n\} \geq \lim_{n \rightarrow \infty} \inf \{a_k : k \geq n\} \equiv \liminf_{n \rightarrow \infty} a_n.$$

Suppose first that $\lim_{n \rightarrow \infty} a_n$ exists and is a real number a . Then from the definition of a limit, there exists N corresponding to $\varepsilon/6$ in the definition. Hence, if $m, n \geq N$, then

$$|a_n - a_m| \leq |a_n - a| + |a - a_m| < \frac{\varepsilon}{6} + \frac{\varepsilon}{6} = \frac{\varepsilon}{3}.$$

From the definition of $\sup\{a_k : k \geq N\}$, there exists $n_1 \geq N$ such that

$$\sup\{a_k : k \geq N\} \leq a_{n_1} + \varepsilon/3.$$

Similarly, there exists $n_2 \geq N$ such that $\inf\{a_k : k \geq N\} \geq a_{n_2} - \varepsilon/3$. It follows that

$$\sup\{a_k : k \geq N\} - \inf\{a_k : k \geq N\} \leq |a_{n_1} - a_{n_2}| + \frac{2\varepsilon}{3} < \varepsilon.$$

Since the sequence, $\{\sup\{a_k : k \geq N\}\}_{N=1}^\infty$ is decreasing and $\{\inf\{a_k : k \geq N\}\}_{N=1}^\infty$ is increasing, it follows that

$$0 \leq \lim_{N \rightarrow \infty} \sup\{a_k : k \geq N\} - \lim_{N \rightarrow \infty} \inf\{a_k : k \geq N\} \leq \varepsilon$$

Since ε is arbitrary, this shows

$$\lim_{N \rightarrow \infty} \sup\{a_k : k \geq N\} = \lim_{N \rightarrow \infty} \inf\{a_k : k \geq N\} \quad (2.1)$$

Next suppose 2.1 and both equal $a \in \mathbb{R}$. Then

$$\lim_{N \rightarrow \infty} (\sup\{a_k : k \geq N\} - \inf\{a_k : k \geq N\}) = 0$$

Since $\sup\{a_k : k \geq N\} \geq \inf\{a_k : k \geq N\}$, it follows that for every $\varepsilon > 0$, there exists N such that $\sup\{a_k : k \geq N\} - \inf\{a_k : k \geq N\} < \varepsilon$, and for every N , $\inf\{a_k : k \geq N\} \leq a \leq \sup\{a_k : k \geq N\}$

$$\inf\{a_k : k \geq N\} \leq a \leq \sup\{a_k : k \geq N\}$$

Thus if $n \geq N$, $|a - a_n| < \varepsilon$ which implies that $\lim_{n \rightarrow \infty} a_n = a$. In case

$$a = \infty = \lim_{N \rightarrow \infty} \sup\{a_k : k \geq N\} = \lim_{N \rightarrow \infty} \inf\{a_k : k \geq N\}$$

then if $r \in \mathbb{R}$ is given, there exists N such that $\inf\{a_k : k \geq N\} > r$ which is to say that $\lim_{n \rightarrow \infty} a_n = \infty$. The case where $a = -\infty$ is similar except you use $\sup\{a_k : k \geq N\}$. ■

The significance of \limsup and \liminf , in addition to what was just discussed, is contained in the following theorem which follows quickly from the definition.

Theorem 2.6.5 Suppose $\{a_n\}$ is a sequence of points of $[-\infty, \infty]$. Also define $\lambda = \limsup_{n \rightarrow \infty} a_n$. Then if $b > \lambda$, it follows there exists N such that whenever $n \geq N$, $a_n \leq b$. If $c < \lambda$, then $a_n > c$ for infinitely many values of n . Let $\gamma = \liminf_{n \rightarrow \infty} a_n$. Then if $d < \gamma$, it follows there exists N such that whenever $n \geq N$, $a_n \geq d$. If $e > \gamma$, it follows $a_n < e$ for infinitely many values of n .

The proof of this theorem is left as an exercise for you. It follows directly from the definition and it is the sort of thing you must do yourself. Here is one other simple proposition.

Proposition 2.6.6 Let $\lim_{n \rightarrow \infty} a_n = a > 0$. Then $\limsup_{n \rightarrow \infty} a_n b_n = a \limsup_{n \rightarrow \infty} b_n$.

Proof: This follows from the definition. Let $\lambda_n = \sup\{a_k b_k : k \geq n\}$. For all n large enough, $a_n > a - \varepsilon$ where ε is small enough that $a - \varepsilon > 0$. Therefore,

$$\lambda_n \geq \sup\{b_k : k \geq n\} (a - \varepsilon)$$

for all n large enough. Then

$$\begin{aligned} \limsup_{n \rightarrow \infty} a_n b_n &= \lim_{n \rightarrow \infty} \lambda_n \equiv \limsup_{n \rightarrow \infty} a_n b_n \geq \lim_{n \rightarrow \infty} (\sup \{b_k : k \geq n\} (a - \varepsilon)) \\ &= (a - \varepsilon) \limsup_{n \rightarrow \infty} b_n \end{aligned}$$

Similar reasoning shows $\limsup_{n \rightarrow \infty} a_n b_n \leq (a + \varepsilon) \limsup_{n \rightarrow \infty} b_n$. Now since $\varepsilon > 0$ is arbitrary, the conclusion follows. ■

2.7 Nested Interval Lemma

The nested interval lemma is a simple and important lemma which is used later quite a bit.

Lemma 2.7.1 *Let $[a_k, b_k] \supseteq [a_{k+1}, b_{k+1}]$ for all $k = 1, 2, 3, \dots$. Then there exists a point p in $\bigcap_{k=1}^{\infty} [a_k, b_k]$. If $\lim_{k \rightarrow \infty} (b_k - a_k) = 0$, then there is only one such point*

Proof: We note that for any $k, l, a_k \leq b_l$. Here is why. If $k \leq l$, then $a_k \leq a_l \leq b_l$. If $k > l$, then $b_l \geq b_k \geq a_k$. It follows that for each l , $\sup_k a_k \leq b_l$. Hence $\sup_k a_k$ is a lower bound to the set of all b_l and so it is no larger than the greatest lower bound. It follows that $\sup_k a_k \leq \inf_l b_l$. Pick $x \in [\sup_k a_k, \inf_l b_l]$. Then for every $k, a_k \leq x \leq b_k$. Hence $x \in \bigcap_{k=1}^{\infty} [a_k, b_k]$.

To see the last claim, if q is another point in all the intervals, then both p and q are in $[a_k, b_k]$ and so $|p - q| \leq (b_k - a_k) < \varepsilon$ if k is large enough. Since ε is arbitrary, $p = q$. ■

2.8 The Hausdorff Maximal Theorem

This major theorem, or something like it (Several equivalent statements are proved later.), is either absolutely essential or extremely convenient. First is the definition of what is meant by a partial order.

Definition 2.8.1 *A nonempty set \mathcal{F} is called a partially ordered set if it has a partial order denoted by \prec . This means it satisfies the following. If $x \prec y$ and $y \prec z$, then $x \prec z$. Also $x \prec x$. It is like \subseteq on the set of all subsets of a given set. It is not the case that given two elements of \mathcal{F} that they are related. In other words, you cannot conclude that either $x \prec y$ or $y \prec x$. A chain, denoted by $\mathcal{C} \subseteq \mathcal{F}$ has the property that it is totally ordered meaning that if $x, y \in \mathcal{C}$, either $x \prec y$ or $y \prec x$. A maximal chain is a chain \mathcal{C} which has the property that there is no strictly larger chain. In other words, if $x \in \mathcal{F} \setminus \mathcal{C}$, then $\mathcal{C} \cup \{x\}$ is no longer a chain.*

Here is the Hausdorff maximal theorem. The proof is a proof by contradiction. We assume there is no maximal chain and then show this cannot happen. The axiom of choice is used in choosing the $x_{\mathcal{C}}$ right at the beginning of the argument.

Theorem 2.8.2 *Let \mathcal{F} be a nonempty partially ordered set with order \prec . Then there exists a maximal chain.*

Proof: Suppose not. Then for \mathcal{C} a chain, let $\theta\mathcal{C}$ denote $\mathcal{C} \cup \{x_{\mathcal{C}}\}$. Thus for \mathcal{C} a chain, $\theta\mathcal{C}$ is a larger chain which has exactly one more element of \mathcal{F} . Since $\mathcal{F} \neq \emptyset$, pick $x_0 \in \mathcal{F}$. Note that $\{x_0\}$ is a chain. Let \mathcal{X} be the set of all chains \mathcal{C} such that $x_0 \in \mathcal{C}$. Thus \mathcal{X} contains $\{x_0\}$. Call two chains comparable if one is a subset of the other. Also, if \mathcal{S}

is a nonempty subset of \mathcal{F} in which all chains are comparable, then $\cup \mathcal{S}$ is also a chain. From now on \mathcal{S} **will always refer to a nonempty set of chains in which any pair are comparable**. Then summarizing,

1. $x_0 \in \cup \mathcal{C}$ for all $\mathcal{C} \in \mathcal{X}$.
2. $\{x_0\} \in \mathcal{X}$
3. If $\mathcal{C} \in \mathcal{X}$ then $\theta \mathcal{C} \in \mathcal{X}$.
4. If $\mathcal{S} \subseteq \mathcal{X}$ then $\cup \mathcal{S} \in \mathcal{X}$.

A subset \mathcal{Y} of \mathcal{X} will be called a “tower” if \mathcal{Y} satisfies 1.) - 4.). Let \mathcal{Y}_0 be the intersection of all towers. Then \mathcal{Y}_0 is also a tower, the smallest one. Then the next claim might seem to be so because if not, \mathcal{Y}_0 would not be the smallest tower.

Claim 1: If $\mathcal{C}_0 \in \mathcal{Y}_0$ is comparable to every chain $\mathcal{C} \in \mathcal{Y}_0$, then if $\mathcal{C}_0 \subsetneq \mathcal{C}$, it must be the case that $\theta \mathcal{C}_0 \subseteq \mathcal{C}$. In other words, $x_{\mathcal{C}_0} \in \cup \mathcal{C}$. The symbol \subsetneq indicates proper subset.

This is done by considering a set $\mathcal{B} \subseteq \mathcal{Y}_0$ consisting of \mathcal{D} which acts like \mathcal{C} in the above and showing that it actually equals \mathcal{Y}_0 because it is a tower.

Proof of Claim 1: Consider $\mathcal{B} \equiv \{\mathcal{D} \in \mathcal{Y}_0 : \mathcal{D} \subseteq \mathcal{C}_0 \text{ or } x_{\mathcal{C}_0} \in \cup \mathcal{D}\}$. Let $\mathcal{Y}_1 \equiv \mathcal{Y}_0 \cap \mathcal{B}$. I want to argue that \mathcal{Y}_1 is a tower. By definition all chains of \mathcal{Y}_1 contain x_0 in their unions. If $\mathcal{D} \in \mathcal{Y}_1$, is $\theta \mathcal{D} \in \mathcal{Y}_1$? If $\mathcal{S} \subseteq \mathcal{Y}_1$, is $\cup \mathcal{S} \in \mathcal{Y}_1$? Is $\{x_0\} \in \mathcal{B}$?

$\{x_0\}$ cannot properly contain \mathcal{C}_0 since $x_0 \in \cup \mathcal{C}_0$. Therefore, $\mathcal{C}_0 \supseteq \{x_0\}$ so $\{x_0\} \in \mathcal{B}$.

If $\mathcal{S} \subseteq \mathcal{Y}_1$, and $\mathcal{D} \equiv \cup \mathcal{S}$, is $\mathcal{D} \in \mathcal{Y}_1$? Since \mathcal{Y}_0 is a tower, \mathcal{D} is comparable to \mathcal{C}_0 . If $\mathcal{D} \subseteq \mathcal{C}_0$, then \mathcal{D} is in \mathcal{B} . Otherwise $\mathcal{D} \supsetneq \mathcal{C}_0$ and in this case, why is \mathcal{D} in \mathcal{B} ? Why is $x_{\mathcal{C}_0} \in \cup \mathcal{D}$? The chains of \mathcal{S} are in \mathcal{B} so one of them, called \mathcal{C} must properly contain \mathcal{C}_0 and so $x_{\mathcal{C}_0} \in \cup \mathcal{C} \subseteq \cup \mathcal{D}$. Therefore, $\mathcal{D} \in \mathcal{B} \cap \mathcal{Y}_0 = \mathcal{Y}_1$. 4.) holds. Two cases remain, to show that \mathcal{Y}_1 satisfies 3.).

case 1: $\mathcal{D} \supsetneq \mathcal{C}_0$. Then by definition of \mathcal{B} , $x_{\mathcal{C}_0} \in \cup \mathcal{D}$ and so $x_{\mathcal{C}_0} \in \cup \theta \mathcal{D}$ so $\theta \mathcal{D} \in \mathcal{Y}_1$.

case 2: $\mathcal{D} \subseteq \mathcal{C}_0$. $\theta \mathcal{D} \in \mathcal{Y}_0$ so $\theta \mathcal{D}$ is comparable to \mathcal{C}_0 . First suppose $\theta \mathcal{D} \supsetneq \mathcal{C}_0$. Thus $\mathcal{D} \subseteq \mathcal{C}_0 \subsetneq \mathcal{D} \cup \{x\}$. If $x \in \mathcal{C}_0$ and x is not in \mathcal{D} then $\mathcal{D} \cup \{x\} \subseteq \mathcal{C}_0 \subsetneq \mathcal{D} \cup \{x_{\mathcal{D}}\}$. This is impossible. Consider x . Thus in this case that $\theta \mathcal{D} \supsetneq \mathcal{C}_0$, $\mathcal{D} = \mathcal{C}_0$. It follows that $x_{\mathcal{D}} = x_{\mathcal{C}_0} \in \cup \theta \mathcal{C}_0 = \cup \theta \mathcal{D}$ and so $\theta \mathcal{D} \in \mathcal{Y}_1$. The other case is that $\theta \mathcal{D} \subseteq \mathcal{C}_0$ so $\theta \mathcal{D} \in \mathcal{B}$ by definition. This shows 3.) so \mathcal{Y}_1 is a tower and must equal \mathcal{Y}_0 .

Claim 2: Any two chains in \mathcal{Y}_0 are comparable.

Proof of Claim 2: Let \mathcal{Y}_1 consist of all chains of \mathcal{Y}_0 which are comparable to every chain of \mathcal{Y}_0 . $\{x_0\}$ is in \mathcal{Y}_1 by definition. All chains of \mathcal{Y}_0 have x_0 in their union. If $\mathcal{S} \subseteq \mathcal{Y}_1$, is $\cup \mathcal{S} \in \mathcal{Y}_1$? Given $\mathcal{D} \in \mathcal{Y}_0$ either every chain of \mathcal{S} is contained in \mathcal{D} or at least one contains \mathcal{D} . Either way \mathcal{D} is comparable to $\cup \mathcal{S}$ so $\cup \mathcal{S} \in \mathcal{Y}_1$. It remains to show 3.). Let $\mathcal{C} \in \mathcal{Y}_1$ and $\mathcal{D} \in \mathcal{Y}_0$. Since \mathcal{C} is comparable to all chains in \mathcal{Y}_0 , it follows from Claim 1 either $\mathcal{C} \subsetneq \mathcal{D}$ when $x_{\mathcal{C}} \in \cup \mathcal{D}$ and $\theta \mathcal{C} \subseteq \mathcal{D}$ or $\mathcal{C} \supsetneq \mathcal{D}$ when $\theta \mathcal{C} \supsetneq \mathcal{D}$. Hence $\mathcal{Y}_1 = \mathcal{Y}_0$ because \mathcal{Y}_0 is as small as possible.

Since every pair of chains in \mathcal{Y}_0 are comparable and \mathcal{Y}_0 is a tower, it follows that $\cup \mathcal{Y}_0 \in \mathcal{Y}_0$ so $\cup \mathcal{Y}_0$ is a chain. However, $\theta \cup \mathcal{Y}_0$ is a chain which properly contains $\cup \mathcal{Y}_0$ and since \mathcal{Y}_0 is a tower, $\theta \cup \mathcal{Y}_0 \in \mathcal{Y}_0$. Thus $\cup (\theta \cup \mathcal{Y}_0) \supsetneq \cup (\cup \mathcal{Y}_0) \supsetneq \cup (\theta \cup \mathcal{Y}_0)$ which is a contradiction. Therefore, for some chain \mathcal{C} it is impossible to obtain the $x_{\mathcal{C}}$ described above and so, this \mathcal{C} is a maximal chain. ■

If X is a nonempty set, \leq is an order on X if

$$\begin{aligned} & x \leq x, \\ & \text{either } x \leq y \text{ or } y \leq x \\ & \text{if } x \leq y \text{ and } y \leq z \text{ then } x \leq z. \end{aligned}$$

and \leq is a well order if (X, \leq) if every nonempty subset of X has a smallest element. More precisely, if $S \neq \emptyset$ and $S \subseteq X$ then there exists an $x \in S$ such that $x \leq y$ for all $y \in S$. A familiar example of a well-ordered set is the natural numbers.

Lemma 2.8.3 *The Hausdorff maximal principle implies every nonempty set can be well-ordered.*

Proof: Let X be a nonempty set and let $a \in X$. Then $\{a\}$ is a well-ordered subset of X . Let $\mathcal{F} = \{S \subseteq X : \text{there exists a well order for } S\}$. Thus $\mathcal{F} \neq \emptyset$. For $S_1, S_2 \in \mathcal{F}$, define $S_1 \prec S_2$ if $S_1 \subseteq S_2$ and there exists a well order for S_2 , \leq_2 such that (S_2, \leq_2) is well-ordered and if $y \in S_2 \setminus S_1$ then $x \leq_2 y$ for all $x \in S_1$, and if \leq_1 is the well order of S_1 then the two orders are consistent on S_1 . Then observe that \prec is a partial order on \mathcal{F} . By the Hausdorff maximal principle, let \mathcal{C} be a maximal chain in \mathcal{F} and let $X_\infty \equiv \cup \mathcal{C}$. Define an order, \leq , on X_∞ as follows. If x, y are elements of X_∞ , pick $S \in \mathcal{C}$ such that x, y are both in S . Then if \leq_S is the order on S , let $x \leq y$ if and only if $x \leq_S y$. This definition is well defined because of the definition of the order, \prec . Now let U be any nonempty subset of X_∞ . Then $S \cap U \neq \emptyset$ for some $S \in \mathcal{C}$. Because of the definition of \leq , if $y \in S_2 \setminus S_1$, $S_i \in \mathcal{C}$, then $x \leq y$ for all $x \in S_1$. Thus, if $y \in X_\infty \setminus S$ then $x \leq y$ for all $x \in S$ and so the smallest element of $S \cap U$ exists and is the smallest element in U . Therefore X_∞ is well-ordered. Now suppose there exists $z \in X \setminus X_\infty$. Define the following order, \leq_1 , on $X_\infty \cup \{z\}$.

$$x \leq_1 y \text{ if and only if } x \leq y \text{ whenever } x, y \in X_\infty$$

$$x \leq_1 z \text{ whenever } x \in X_\infty.$$

Let $\tilde{\mathcal{C}} = \{S \in \mathcal{C} \text{ or } X_\infty \cup \{z\}\}$. Then $\tilde{\mathcal{C}}$ is a strictly larger chain than \mathcal{C} contradicting maximality of \mathcal{C} . Thus $X \setminus X_\infty = \emptyset$ and this shows X is well-ordered by \leq . ■

With these two lemmas the main result follows.

Theorem 2.8.4 *The following are equivalent.*

The axiom of choice

The Hausdorff maximal principle

The well-ordering principle.

Proof: It remains to show that the well-ordering principle implies the axiom of choice. Let I be a nonempty set and let X_i be a nonempty set for each $i \in I$. Let $X = \cup \{X_i : i \in I\}$ and well order X . Let $f(i)$ be the smallest element of X_i . Then $f \in \prod_{i \in I} X_i$. ■

The book by Hewitt and Stromberg [26] has more equivalences.

Chapter 3

Metric Spaces

3.1 Open and Closed Sets, Sequences, Limit Points

It is most efficient to discuss things in terms of abstract metric spaces to begin with.

Definition 3.1.1 A non empty set X is called a metric space if there is a function $d : X \times X \rightarrow [0, \infty)$ which satisfies the following axioms.

1. $d(x, y) = d(y, x)$
2. $d(x, y) \geq 0$ and equals 0 if and only if $x = y$
3. $d(x, y) + d(y, z) \geq d(x, z)$

This function d is called the metric. We often refer to it as the distance also.

Definition 3.1.2 An open ball, denoted as $B(x, r)$ is defined as follows.

$$B(x, r) \equiv \{y : d(x, y) < r\}$$

A set U is said to be open if whenever $x \in U$, it follows that there is $r > 0$ such that $B(x, r) \subseteq U$. More generally, a point x is said to be an interior point of U if there exists such a ball. In words, an open set is one for which every point is an interior point.

For example, you could have X be a subset of \mathbb{R} and $d(x, y) = |x - y|$. Then the first thing to show is the following.

Proposition 3.1.3 An open ball is an open set.

Proof: Suppose $y \in B(x, r)$. We need to verify that y is an interior point of $B(x, r)$. Let $\delta = r - d(x, y)$. Then if $z \in B(y, \delta)$, it follows that

$$d(z, x) \leq d(z, y) + d(y, x) < \delta + d(y, x) = r - d(x, y) + d(y, x) = r$$

Thus $y \in B(y, \delta) \subseteq B(x, r)$. ■

Definition 3.1.4 Let S be a nonempty subset of a metric space. Then p is a limit point (accumulation point) of S if for every $r > 0$ there exists a point different than p in $B(p, r) \cap S$. Sometimes people denote the set of limit points as S' .

The following proposition is fairly obvious from the above definition and will be used whenever convenient. It is equivalent to the above definition and so it can take the place of the above definition if desired.

Proposition 3.1.5 A point x is a limit point of the nonempty set A if and only if every $B(x, r)$ contains infinitely many points of A .

Proof: \Leftarrow is obvious. Consider \Rightarrow . Let x be a limit point. Let $r_1 = 1$. Then $B(x, r_1)$ contains $a_1 \neq x$. If $\{a_1, \dots, a_n\}$ have been chosen none equal to x and with no repeats in the list, let $0 < r_n < \min(\frac{1}{n}, \min\{d(a_i, x), i = 1, 2, \dots, n\})$. Then let $a_{n+1} \in B(x, r_n)$. Thus every $B(x, r)$ contains $B(x, r_n)$ for all n large enough and hence it contains a_k for $k \geq n$ where the a_k are distinct, none equal to x . ■

A related idea is the notion of the limit of a sequence. Recall that a sequence is really just a mapping from \mathbb{N} to X . We write them as $\{x_n\}$ or $\{x_n\}_{n=1}^\infty$ if we want to emphasize the values of n . Then the following definition is what it means for a sequence to converge.

Definition 3.1.6 We say that $x = \lim_{n \rightarrow \infty} x_n$ when for every $\varepsilon > 0$ there exists N such that if $n \geq N$, then

$$d(x, x_n) < \varepsilon$$

Often we write $x_n \rightarrow x$ for short. This is equivalent to saying

$$\lim_{n \rightarrow \infty} d(x, x_n) = 0.$$

Proposition 3.1.7 The limit is well defined. That is, if x, x' are both limits of a sequence, then $x = x'$.

Proof: From the definition, there exist N, N' such that if $n \geq N$, then $d(x, x_n) < \varepsilon/2$ and if $n \geq N'$, then $d(x, x_n) < \varepsilon/2$. Then let $M \geq \max(N, N')$. Let $n > M$. Then

$$d(x, x') \leq d(x, x_n) + d(x_n, x') < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon$$

Since ε is arbitrary, this shows that $x = x'$ because $d(x, x') = 0$. ■

Next there is an important theorem about limit points and convergent sequences.

Theorem 3.1.8 Let $S \neq \emptyset$. Then p is a limit point of S if and only if there exists a sequence of distinct points of $S, \{x_n\}$ none of which equal p such that $\lim_{n \rightarrow \infty} x_n = p$.

Proof: \Rightarrow Suppose p is a limit point. Why does there exist the promised convergent sequence? Let $x_1 \in B(p, 1) \cap S$ such that $x_1 \neq p$. If x_1, \dots, x_n have been chosen, let $x_{n+1} \neq p$ be in $B(p, \delta_{n+1}) \cap S$ where

$$\delta_{n+1} = \min \left\{ \frac{1}{n+1}, d(x_i, p), i = 1, 2, \dots, n \right\}.$$

Then this constructs the necessary convergent sequence.

\Leftarrow Conversely, if such a sequence $\{x_n\}$ exists, then for every $r > 0$, $B(p, r)$ contains $x_n \in S$ for all n large enough. Hence, p is a limit point because none of these x_n are equal to p . ■

Definition 3.1.9 A set H is closed means H^C is open.

Note that this says that the complement of an open set is closed. If V is open, then the complement of its complement is itself. Thus $(V^C)^C = V$ an open set. Hence V^C is closed.

Then the following theorem gives the relationship between closed sets and limit points.

Theorem 3.1.10 A set H is closed if and only if it contains all of its limit points.

Proof: \Rightarrow Let H be closed and let p be a limit point. We need to verify that $p \in H$. If it is not, then since H is closed, its complement is open and so there exists $\delta > 0$ such that $B(p, \delta) \cap H = \emptyset$. However, this prevents p from being a limit point.

\Leftarrow Next suppose H has all of its limit points. Why is H^C open? If $p \in H^C$ then it is not a limit point and so there exists $\delta > 0$ such that $B(p, \delta)$ has no points of H . In other words, H^C is open. Hence H is closed. ■

Corollary 3.1.11 *A set H is closed if and only if whenever $\{h_n\}$ is a sequence of points of H which converges to a point x , it follows that $x \in H$.*

Proof: \Rightarrow Suppose H is closed and $h_n \rightarrow x$. If $x \in H$ there is nothing left to show. If $x \notin H$, then from the definition of limit, it is a limit point of H because none of the h_n are equal to x . Hence $x \in H$ after all.

\Leftarrow Suppose the limit condition holds, why is H closed? Let $x \in H'$ the set of limit points of H . By Theorem 3.1.8 there exists a sequence of points of H , $\{h_n\}$ such that $h_n \rightarrow x$. Then by assumption, $x \in H$. Thus H contains all of its limit points and so it is closed by Theorem 3.1.10. ■

Next is the important concept of a subsequence.

Definition 3.1.12 *Let $\{x_n\}_{n=1}^\infty$ be a sequence. Then if $n_1 < n_2 < \dots$ is a strictly increasing sequence of indices, we say $\{x_{n_k}\}_{k=1}^\infty$ is a subsequence of $\{x_n\}_{n=1}^\infty$.*

The really important thing about subsequences is that they preserve convergence.

Theorem 3.1.13 *Let $\{x_{n_k}\}$ be a subsequence of a convergent sequence $\{x_n\}$ where $x_n \rightarrow x$. Then $\lim_{k \rightarrow \infty} x_{n_k} = x$ also.*

Proof: Let $\varepsilon > 0$ be given. Then there exists N such that $d(x_n, x) < \varepsilon$ if $n \geq N$. It follows that if $k \geq N$, then $n_k \geq N$ and so $d(x_{n_k}, x) < \varepsilon$ if $k \geq N$. This is what it means to say $\lim_{k \rightarrow \infty} x_{n_k} = x$. ■

3.2 Cauchy Sequences, Completeness

Of course it does not go the other way. For example, you could let $x_n = (-1)^n$ and it has a convergent subsequence but fails to converge. Here $d(x, y) = |x - y|$ and the metric space is just \mathbb{R} .

However, there is a kind of sequence for which it does go the other way. This is called a Cauchy sequence.

Definition 3.2.1 *$\{x_n\}$ is called a Cauchy sequence if for every $\varepsilon > 0$ there exists N such that if $m, n \geq N$, then $d(x_n, x_m) < \varepsilon$.*

Now the major theorem about this is the following.

Theorem 3.2.2 *Let $\{x_n\}$ be a Cauchy sequence. Then it converges if and only if any subsequence converges.*

Proof: \Rightarrow This was just done above. \Leftarrow Suppose now that $\{x_n\}$ is a Cauchy sequence and $\lim_{k \rightarrow \infty} x_{n_k} = x$. Then there exists N_1 such that if $k > N_1$, then $d(x_{n_k}, x) < \varepsilon/2$. From the definition of what it means to be Cauchy, there exists N_2 such that if $m, n \geq N_2$, then $d(x_m, x_n) < \varepsilon/2$. Let $N \geq \max(N_1, N_2)$. Then if $k \geq N$, then $n_k \geq N$ and so $d(x, x_k) \leq d(x, x_{n_k}) + d(x_{n_k}, x_k) < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon$. It follows from the definition that $\lim_{k \rightarrow \infty} x_k = x$. ■

Definition 3.2.3 A metric space is said to be **complete** if every Cauchy sequence converges.

There certainly are metric spaces which are not complete. For example, if you consider \mathbb{Q} with $d(x, y) \equiv |x - y|$, this will not be complete because you can get a sequence which is obtained as x_n defined as the n decimal place description of $\sqrt{2}$. However, if a sequence converges, then it must be Cauchy.

Lemma 3.2.4 If $x_n \rightarrow x$, then $\{x_n\}$ is a Cauchy sequence.

Proof: Let $\varepsilon > 0$. Then there exists n_ε such that if $m \geq n_\varepsilon$, then $d(x, x_m) < \varepsilon/2$. If $m, k \geq n_\varepsilon$, then by the triangle inequality, $d(x_m, x_k) \leq d(x_m, x) + d(x, x_k) < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon$ showing that the convergent sequence is indeed a Cauchy sequence as claimed. ■

Another nice thing to note is this.

Proposition 3.2.5 If $\{x_n\}$ is a sequence and if p is a limit point of the set $S = \bigcup_{n=1}^{\infty} \{x_n\}$, then there is a subsequence $\{x_{n_k}\}$ such that $\lim_{k \rightarrow \infty} x_{n_k} = p$.

Proof: By Theorem 3.1.8, there exists a sequence of distinct points of S denoted as $\{y_k\}$ such that none of them equal p and $\lim_{k \rightarrow \infty} y_k = p$. Thus $B(p, r)$ contains infinitely many different points of the set D , this for every r . Let $x_{n_1} \in B(p, 1)$ where n_1 is the first index such that $x_{n_1} \in B(p, 1)$. Suppose x_{n_1}, \dots, x_{n_k} have been chosen, the n_i increasing and let $1 > \delta_1 > \delta_2 > \dots > \delta_k$ where $x_{n_i} \in B(p, \delta_i)$. Then let

$$\delta_{k+1} < \min \left\{ \frac{1}{2^{k+1}}, d(p, x_{n_j}), \delta_j, j = 1, 2, \dots, k \right\}$$

Let $x_{n_{k+1}} \in B(p, \delta_{k+1})$ where n_{k+1} is the first index such that $x_{n_{k+1}}$ is contained $B(p, \delta_{k+1})$. Then $\lim_{k \rightarrow \infty} x_{n_k} = p$. ■

Another useful result is the following.

Lemma 3.2.6 Suppose $x_n \rightarrow x$ and $y_n \rightarrow y$. Then $d(x_n, y_n) \rightarrow d(x, y)$.

Proof: Consider the following.

$$d(x, y) \leq d(x, x_n) + d(x_n, y) \leq d(x, x_n) + d(x_n, y_n) + d(y_n, y)$$

so $d(x, y) - d(x_n, y_n) \leq d(x, x_n) + d(y_n, y)$. Similar reasoning to what was just used shows that $d(x_n, y_n) - d(x, y) \leq d(x, x_n) + d(y_n, y)$, so $|d(x_n, y_n) - d(x, y)| \leq d(x, x_n) + d(y_n, y)$ and the right side converges to 0 as $n \rightarrow \infty$. ■

3.3 Closure of a Set

Next is the topic of the closure of a set.

Definition 3.3.1 Let A be a nonempty subset of (X, d) a metric space. Then \bar{A} is defined to be the intersection of all closed sets which contain A . Note the whole space, X is one such closed set which contains A . The whole space X is closed because its complement is open, its complement being \emptyset . It is certainly true that every point of the empty set is an interior point because there are no points of \emptyset .

Lemma 3.3.2 *Let A be a nonempty set in (X, d) . Then \bar{A} is a closed set and $\bar{A} = A \cup A'$ where A' denotes the set of limit points of A .*

Proof: First of all, denote by \mathcal{C} the set of closed sets which contain A . Then $\bar{A} = \cap \mathcal{C}$ and this will be closed if its complement is open. However, $\bar{A}^C = \cup \{H^C : H \in \mathcal{C}\}$. Each H^C is open and so the union of all these open sets must also be open. This is because if x is in this union, then it is in at least one of them. Hence it is an interior point of that one. But this implies it is an interior point of the union of them all which is an even larger set. Thus \bar{A} is closed.

The interesting part is the next claim. First note that from the definition, $A \subseteq \bar{A}$ so if $x \in A$, then $x \in \bar{A}$. Now consider $y \in A'$ but $y \notin A$. If $y \notin \bar{A}$, a closed set, then there exists $B(y, r) \subseteq \bar{A}^C$. Thus y cannot be a limit point of A , a contradiction. Therefore, $A \cup A' \subseteq \bar{A}$.

Next suppose $x \in \bar{A}$ and suppose $x \notin A$. Then if $B(x, r)$ contains no points of A different than x , since x itself is not in A , it would follow that $B(x, r) \cap A = \emptyset$ and so recalling that open balls are open, $B(x, r)^C$ is a closed set containing A so from the definition, it also contains \bar{A} which is contrary to the assertion that $x \in \bar{A}$. Hence if $x \notin A$, then $x \in A'$ and so $A \cup A' \subseteq \bar{A}$ ■

3.4 Separable Metric Spaces

Definition 3.4.1 *A metric space is called separable if there exists a countable dense subset D . This means two things. First, D is countable, and second, that if x is any point and $r > 0$, then $B(x, r) \cap D \neq \emptyset$. A metric space is called completely separable if there exists a countable collection of nonempty open sets \mathcal{B} such that every open set is the union of some subset of \mathcal{B} . This collection of open sets is called a countable basis.*

For those who like to fuss about empty sets, the empty set is open and it is indeed the union of a subset of \mathcal{B} namely the empty subset.

Theorem 3.4.2 *A metric space is separable if and only if it is completely separable.*

Proof: \Leftarrow Let \mathcal{B} be the special countable collection of open sets and for each $B \in \mathcal{B}$, let p_B be a point of B . Then let $\mathcal{P} \equiv \{p_B : B \in \mathcal{B}\}$. If $B(x, r)$ is any ball, then it is the union of sets of \mathcal{B} and so there is a point of \mathcal{P} in it. Since \mathcal{B} is countable, so is \mathcal{P} .

\Rightarrow Let D be the countable dense set and let $\mathcal{B} \equiv \{B(d, r) : d \in D, r \in \mathbb{Q} \cap [0, \infty)\}$. Then \mathcal{B} is countable because the Cartesian product of countable sets is countable. It suffices to show that every ball is the union of these sets. Let $B(x, R)$ be a ball. Let $y \in B(y, \delta) \subseteq B(x, R)$. Then there exists $d \in B\left(y, \frac{\delta}{10}\right)$. Let $\varepsilon \in \mathbb{Q}$ and $\frac{\delta}{10} < \varepsilon < \frac{\delta}{5}$. Then $y \in B(d, \varepsilon) \subseteq \mathcal{B}$. Is $B(d, \varepsilon) \subseteq B(x, R)$? If so, then the desired result follows because this would show that every $y \in B(x, R)$ is contained in one of these sets of \mathcal{B} which is contained in $B(x, R)$ showing that $B(x, R)$ is the union of sets of \mathcal{B} . Let $z \in B(d, \varepsilon) \subseteq B\left(d, \frac{\delta}{5}\right)$. Then

$$d(y, z) \leq d(y, d) + d(d, z) < \frac{\delta}{10} + \varepsilon < \frac{\delta}{10} + \frac{\delta}{5} < \delta$$

Hence $B(d, \varepsilon) \subseteq B(y, \delta) \subseteq B(x, R)$. Therefore, every ball is the union of sets of \mathcal{B} and, since every open set is the union of balls, it follows that every open set is the union of sets of \mathcal{B} . ■

Corollary 3.4.3 *If (X, d) is a metric space and S is a nonempty subset of X , then S is also separable.*

Proof: Let \mathcal{B} be a countable basis for (X, d) . Say \mathcal{B}_S be those sets of \mathcal{B} which have nonempty intersections with S . By axiom of choice, there is a point in each of these intersections. The resulting countable selection of points must be dense in S . Indeed, if $x \in S$, then $B(x, r)$ is the union of sets of \mathcal{B} and so some point just described is in $B(x, r)$. ■

Definition 3.4.4 *Let S be a nonempty set. Then a set of open sets \mathcal{C} is called an open cover of S if $\bigcup \mathcal{C} \supseteq S$. (It covers up the set S . Think lilly pads covering the surface of a pond.)*

One of the important properties possessed by separable metric spaces is the Lindeloff property.

Definition 3.4.5 *A metric space has the Lindeloff property if whenever \mathcal{C} is an open cover of a set S , there exists a countable subset of \mathcal{C} denoted here by \mathcal{B} such that \mathcal{B} is also an open cover of S .*

Theorem 3.4.6 *Every separable metric space has the Lindeloff property.*

Proof: Let \mathcal{C} be an open cover of a set S . Let \mathcal{B} be a countable basis. Such exists by Theorem 3.4.2. Let \mathcal{B} denote those sets of \mathcal{B} which are contained in some set of \mathcal{C} . Thus \mathcal{B} is a countable open cover of S . Now for $B \in \mathcal{B}$, let U_B be a set of \mathcal{C} which contains B . Letting $\hat{\mathcal{C}}$ denote these sets U_B it follows that $\hat{\mathcal{C}}$ is countable and is an open cover of S . ■

Definition 3.4.7 *A Polish space is a complete separable metric space. These things turn out to be very useful in probability theory and in other areas.*

3.5 Compact Sets

As usual, we are not worrying about empty sets. Fussing over these is usually a waste of time. Thus if a set is mentioned, the default is that it is nonempty.

Definition 3.5.1 *A metric space K is compact if whenever \mathcal{C} is an open cover of K , meaning $K \subseteq \bigcup \mathcal{C}$, there exists a finite subset of \mathcal{C} $\{U_1, \dots, U_n\}$ such that $K \subseteq \bigcup_{k=1}^n U_k$. In words, every open cover admits a finite sub-cover.*

Directly from this definition is the following proposition.

Proposition 3.5.2 *If K is a closed, nonempty subset of a nonempty compact set H , then K is compact.*

Proof: Let \mathcal{C} be an open cover for K . Then $\mathcal{C} \cup \{K^C\}$ is an open cover for H . Thus there are finitely many sets from this last collection of open sets, U_1, \dots, U_m which covers H . Include only those which are in \mathcal{C} . These cover K because K^C covers no points of K . ■

This is the real definition given above. However, in metric spaces, it is equivalent to another definition called sequentially compact.

Definition 3.5.3 A metric space K is sequentially compact means that whenever $\{x_n\} \subseteq K$, there exists a subsequence $\{x_{n_k}\}$ such that $\lim_{k \rightarrow \infty} x_{n_k} = x \in K$ for some point x . In words, every sequence has a subsequence which converges to a point in the set.

There is a fundamental property possessed by a sequentially compact set in a metric space which is described in the following proposition. The special number described is called a Lebesgue number.

Proposition 3.5.4 Let K be a sequentially compact set in a metric space and let \mathcal{C} be an open cover of K . Then there exists a number $\delta > 0$ such that whenever $x \in K$, it follows that $B(x, \delta)$ is contained in some set of \mathcal{C} .

Proof: If \mathcal{C} is an open cover of K and has no Lebesgue number, then for each $n \in \mathbb{N}$, $\frac{1}{n}$ is not a Lebesgue number. Hence there exists $x_n \in K$ such that $B(x_n, \frac{1}{n})$ is not contained in any set of \mathcal{C} . By sequential compactness, there is a subsequence $\{x_{n_k}\}$ such that $x_{n_k} \rightarrow x \in K$. Now there is $r > 0$ such that $B(x, r) \subseteq U \in \mathcal{C}$. Let k be large enough that $\frac{1}{n_k} < \frac{r}{2}$ and also large enough that $x_{n_k} \in B(x, \frac{r}{2})$. Then $B(x_{n_k}, \frac{1}{n_k}) \subseteq B(x_{n_k}, \frac{r}{2}) \subseteq B(x, r)$ contrary to the requirement that $B(x_{n_k}, \frac{1}{n_k})$ is not contained in any set of \mathcal{C} . ■

In any metric space, these two definitions of compactness are equivalent.

Theorem 3.5.5 Let K be a nonempty subset of a metric space (X, d) . Then it is compact if and only if it is sequentially compact.

Proof: \Leftarrow Suppose K is sequentially compact. Let \mathcal{C} be an open cover of K . By Proposition 3.5.4 there is a Lebesgue number $\delta > 0$. Let $x_1 \in K$. If $B(x_1, \delta)$ covers K , then pick a set of \mathcal{C} containing this ball and this set will be a finite subset of \mathcal{C} which covers K . If $B(x_1, \delta)$ does not cover K , let $x_2 \notin B(x_1, \delta)$. Continue this way obtaining x_k such that $d(x_k, x_j) \geq \delta$ whenever $k \neq j$. Thus eventually $\{B(x_i, \delta)\}_{i=1}^n$ must cover K because if not, you could get a sequence $\{x_k\}$ which has every pair of points further apart than δ and hence it has no Cauchy subsequence. Therefore, by Lemma 3.2.4, it would have no convergent subsequence. This would contradict K is sequentially compact. Now let $U_i \in \mathcal{C}$ with $U_i \supseteq B(x_i, \delta)$. Then $\cup_{i=1}^n U_i \supseteq K$.

\Rightarrow Now suppose K is compact. If it is not sequentially compact, then there exists a sequence $\{x_n\}$ which has no convergent subsequence to a point of K . In particular, no point of this sequence is repeated infinitely often. By Proposition 3.2.5 the set of points $\cup_n \{x_n\}$ has no limit point in K . (If it did, you would have a subsequence converging to this point since every ball containing this point would contain infinitely many points of $\cup_n \{x_n\}$.) Now consider the sets $H_n \equiv \cup_{k \geq n} \{x_k\} \cup H'$ where H' denotes all limit points of $\cup_n \{x_n\}$ in X which is the same as the limit points of $\cup_{k \geq n} \{x_k\}$. Therefore, each H_n is closed thanks to Lemma 3.3.2. Now let $U_n \equiv H_n^C$. This is an increasing sequence of open sets whose union contains K thanks to the fact that there is no constant subsequence. However, none of these open sets covers K because U_n is missing x_n , violating the definition of compactness. Next is an alternate argument.

\Rightarrow Now suppose K is compact. If it is not sequentially compact, then there exists a sequence $\{x_n\}$ which has no convergent subsequence to a point of K . If $x \in K$, then there exists $B(x, r_x)$ which contains x_n for only finitely many n . This is because x is not the limit of a subsequence. Then $\{B(x_i, r_i)\}_{i=1}^N$ is a finite sub-cover of K . If p is the largest index for any x_k contained in $\cup_{i=1}^N B(x_i, r_i)$, let $n > p$ and consider x_n . It is a point in K but it can't be in any of the sets covering K . ■

Definition 3.5.6 *X be a metric space. Then a finite set of points $\{x_1, \dots, x_n\}$ is called an ε net if $X \subseteq \bigcup_{k=1}^n B(x_k, \varepsilon)$. If, for every $\varepsilon > 0$ a metric space has an ε net, then we say that the metric space is totally bounded.*

Lemma 3.5.7 *If a metric space (K, d) is sequentially compact, then it is separable and totally bounded.*

Proof: Pick $x_1 \in K$. If $B(x_1, \varepsilon) \supseteq K$, then stop. Otherwise, pick $x_2 \notin B(x_1, \varepsilon)$. Continue this way. If $\{x_1, \dots, x_n\}$ have been chosen, either $K \subseteq \bigcup_{k=1}^n B(x_k, \varepsilon)$ in which case, you have found an ε net or this does not happen in which case, you can pick $x_{n+1} \notin \bigcup_{k=1}^n B(x_k, \varepsilon)$. The process must terminate since otherwise, the sequence would need to have a convergent subsequence which is not possible because every pair of terms is farther apart than ε . See Lemma 3.2.4. Thus for every $\varepsilon > 0$, there is an ε net. Thus the metric space is totally bounded. Let N_ε denote an ε net. Let $D = \bigcup_{k=1}^\infty N_{1/2^k}$. Then this is a countable dense set. It is countable because it is the countable union of finite sets and it is dense because given a point, there is a point of D within $1/2^k$ of it. ■

Also recall that a complete metric space is one for which every Cauchy sequence converges to a point in the metric space.

The following is the main theorem which relates these concepts.

Theorem 3.5.8 *For (X, d) a metric space, the following are equivalent.*

1. (X, d) is compact.
2. (X, d) is sequentially compact.
3. (X, d) is complete and totally bounded.

Proof: By Theorem 3.5.5, the first two conditions are equivalent.

2. \implies 3. If (X, d) is sequentially compact, then by Lemma 3.5.7, it is totally bounded. If $\{x_n\}$ is a Cauchy sequence, then there is a subsequence which converges to $x \in X$ by assumption. However, from Theorem 3.2.2 this requires the original Cauchy sequence to converge.

3. \implies 1. Since (X, d) is totally bounded, there must be a countable dense subset of X . Just take the union of $1/2^k$ nets for each $k \in \mathbb{N}$. Thus (X, d) is completely separable by Theorem 3.4.6 has the Lindeloff property. Hence, if X is not compact, there is a countable set of open sets $\{U_i\}_{i=1}^\infty$ which covers X but no finite subset does. Consider the nonempty closed sets F_n and pick $x_n \in F_n$ where

$$X \setminus \bigcup_{i=1}^n U_i \equiv X \cap (\bigcup_{i=1}^n U_i)^C \equiv F_n$$

Let $\{x_m^k\}_{m=1}^{M_k}$ be a $1/2^k$ net for X . We have for some m , $B(x_m^k, 1/2^k)$ contains x_n for infinitely many values of n because there are only finitely many balls and infinitely many indices. Then out of the finitely many $\{x_m^{k+1}\}$ where $B(x_m^{k+1}, 1/2^{k+1})$ has nonempty intersection with $B(x_m^k, 1/2^k)$, pick one $x_{m_{k+1}}^{k+1}$ such that $B(x_{m_{k+1}}^{k+1}, 1/2^{k+1})$ contains x_n for infinitely many n . Then obviously $\{x_{m_k}^k\}_{k=1}^\infty$ is a Cauchy sequence because

$$d(x_{m_k}^k, x_{m_{k+1}}^{k+1}) \leq \frac{1}{2^k} + \frac{1}{2^{k+1}} \leq \frac{1}{2^{k-1}}$$

Hence for $p < q$,

$$d(x_{m_p}^p, x_{m_q}^q) \leq \sum_{k=p}^{q-1} d(x_{m_k}^k, x_{m_{k+1}}^{k+1}) < \sum_{k=p}^{\infty} \frac{1}{2^{k-1}} = \frac{1}{2^{p-2}}$$

Now take a subsequence $x_{n_k} \in B(x_{m_k}^k, 2^{-k})$ so it follows that $\lim_{k \rightarrow \infty} x_{n_k} = \lim_{k \rightarrow \infty} x_{m_k}^k = x \in X$. However, $x \in F_n$ for each n since each F_n is closed and these sets are nested. Thus $x \in \cap_n F_n$ contrary to the claim that $\{U_i\}_{i=1}^{\infty}$ covers X . ■

For the sake of another point of view, here is another argument, this time that 3.) \Rightarrow 2.). This will illustrate something called the Cantor diagonalization process.

Assume 3.). Suppose $\{x_k\}$ is a sequence in X . By assumption there are finitely many open balls of radius $1/n$ covering X . This for each $n \in \mathbb{N}$. Therefore, for $n = 1$, there is one of the balls, having radius 1 which contains x_k for infinitely many k . Therefore, there is a subsequence with every term contained in this ball of radius 1. Now do for this subsequence what was just done for $\{x_k\}$. There is a further subsequence contained in a ball of radius $1/2$. Continue this way. Denote the i^{th} subsequence as $\{x_{ki}\}_{k=1}^{\infty}$. Arrange them as shown

$$\begin{array}{l} x_{11}, x_{21}, x_{31}, x_{41} \cdots \\ x_{12}, x_{22}, x_{32}, x_{42} \cdots \\ x_{13}, x_{23}, x_{33}, x_{43} \cdots \\ \vdots \end{array}$$

Thus all terms of $\{x_{ki}\}_{k=1}^{\infty}$ are contained in a ball of radius $1/i$. Consider now the diagonal sequence defined as $y_k \equiv x_{kk}$. Given n , each y_k is contained in a ball of radius $1/n$ whenever $k \geq n$. Thus $\{y_k\}$ is a subsequence of the original sequence and $\{y_k\}$ is a Cauchy sequence. By completeness of X , this converges to some $x \in X$ which shows that every sequence in X has a convergent subsequence. This shows 3.) \Rightarrow 2.). ■

Lemma 3.5.9 *The closed interval $[a, b]$ in \mathbb{R} is compact and every Cauchy sequence in \mathbb{R} converges.*

Proof: To show this, suppose it is not. Then there is an open cover \mathcal{C} which admits no finite subcover for $[a, b] \equiv I_0$. Consider the two intervals $[a, \frac{a+b}{2}]$, $[\frac{a+b}{2}, b]$. One of these, maybe both cannot be covered with finitely many sets of \mathcal{C} since otherwise, there would be a finite collection of sets from \mathcal{C} covering $[a, b]$. Let I_1 be the interval which has no finite subcover. Now do for it what was done for I_0 . Split it in half and pick the half which has no finite covering of sets of \mathcal{C} . Thus there is a “nested” sequence of closed intervals $I_0 \supseteq I_1 \supseteq I_2 \cdots$, each being half of the preceding interval. Say $I_n = [a_n, b_n]$. By the nested interval Lemma, Lemma 2.7.1, there is a point x in all these intervals. The point is unique because the lengths of the intervals converge to 0. This point is in some $O \in \mathcal{C}$. Thus for some $\delta > 0$, $[x - \delta, x + \delta]$, having length 2δ , is contained in O . For k large enough, the interval $[a_k, b_k]$ has length less than δ but contains x . Therefore, it is contained in $[x - \delta, x + \delta]$ and so must be contained in a single set of \mathcal{C} contrary to the construction. This contradiction shows that in fact $[a, b]$ is compact.

Now if $\{x_n\}$ is a Cauchy sequence, then it is contained in some interval $[a, b]$ which is compact. Hence there is a subsequence which converges to some $x \in [a, b]$. By Theorem 3.2.2 the original Cauchy sequence converges to x . ■

3.6 Continuous Functions

The following is a fairly general definition of what it means for a function to be continuous. It includes everything seen in typical calculus classes as a special case.

Definition 3.6.1 Let $f : X \rightarrow Y$ be a function where (X, d) and (Y, ρ) are metric spaces. Then f is continuous at $x \in X$ if and only if the following condition holds. For every $\varepsilon > 0$, there exists $\delta > 0$ such that if $d(\hat{x}, x) < \delta$, then $\rho(f(\hat{x}), f(x)) < \varepsilon$. If f is continuous at every $x \in X$ we say that f is continuous on X .

For example, you could have a real valued function $f(x)$ defined on an interval $[0, 1]$. In this case you would have $X = [0, 1]$ and $Y = \mathbb{R}$ with the distance given by $d(x, y) = |x - y|$. Then the following theorem is the main result.

Theorem 3.6.2 Let $f : X \rightarrow Y$ where (X, d) and (Y, ρ) are metric spaces. Then the following two are equivalent.

- a f is continuous at x .
 - b Whenever $x_n \rightarrow x$, it follows that $f(x_n) \rightarrow f(x)$.
- Also, the following are equivalent.
- c f is continuous on X .
 - d Whenever V is open in Y , it follows that $f^{-1}(V) \equiv \{x : f(x) \in V\}$ is open in X .
 - e Whenever H is closed in Y , it follows that $f^{-1}(H) \equiv \{x : f(x) \in H\}$ is closed in X .

Proof: a \implies b: Let f be continuous at x and suppose $x_n \rightarrow x$. Then let $\varepsilon > 0$ be given. By continuity, there exists $\delta > 0$ such that if $d(\hat{x}, x) < \delta$, then $\rho(f(\hat{x}), f(x)) < \varepsilon$. Since $x_n \rightarrow x$, it follows that there exists N such that if $n \geq N$, then $d(x_n, x) < \delta$ and so, if $n \geq N$, it follows that $\rho(f(x_n), f(x)) < \varepsilon$. Since $\varepsilon > 0$ is arbitrary, it follows that $f(x_n) \rightarrow f(x)$.

b \implies a: Suppose b holds but f fails to be continuous at x . Then there exists $\varepsilon > 0$ such that for all $\delta > 0$, there exists \hat{x} such that $d(\hat{x}, x) < \delta$ but $\rho(f(\hat{x}), f(x)) \geq \varepsilon$. Letting $\delta = 1/n$, there exists x_n such that $d(x_n, x) < 1/n$ but $\rho(f(x_n), f(x)) \geq \varepsilon$. Now this is a contradiction because by assumption, the fact that $x_n \rightarrow x$ implies that $f(x_n) \rightarrow f(x)$. In particular, for large enough n , $\rho(f(x_n), f(x)) < \varepsilon$ contrary to the construction.

c \implies d: Let V be open in Y . Let $x \in f^{-1}(V)$ so that $f(x) \in V$. Since V is open, there exists $\varepsilon > 0$ such that $B(f(x), \varepsilon) \subseteq V$. Since f is continuous at x , it follows that there exists $\delta > 0$ such that if $\hat{x} \in B(x, \delta)$, then $f(\hat{x}) \in B(f(x), \varepsilon) \subseteq V$. ($f(B(x, \delta)) \subseteq B(f(x), \varepsilon)$) In other words, $B(x, \delta) \subseteq f^{-1}(B(f(x), \varepsilon)) \subseteq f^{-1}(V)$ which shows that, since x was an arbitrary point of $f^{-1}(V)$, every point of $f^{-1}(V)$ is an interior point which implies $f^{-1}(V)$ is open.

d \implies e: Let H be closed in Y . Then $f^{-1}(H)^C = f^{-1}(H^C)$ which is open by assumption. Hence $f^{-1}(H)$ is closed because its complement is open.

e \implies d: Let V be open in Y . Then $f^{-1}(V)^C = f^{-1}(V^C)$ which is assumed to be closed. This is because the complement of an open set is a closed set.

d \implies c: Let $x \in X$ be arbitrary. Is it the case that f is continuous at x ? Let $\varepsilon > 0$ be given. Then $B(f(x), \varepsilon)$ is an open set in Y and so $x \in f^{-1}(B(f(x), \varepsilon))$ which is given to be open. Hence there exists $\delta > 0$ such that $x \in B(x, \delta) \subseteq f^{-1}(B(f(x), \varepsilon))$. Thus, $f(B(x, \delta)) \subseteq B(f(x), \varepsilon)$ so $\rho(f(\hat{x}), f(x)) < \varepsilon$. Thus f is continuous at x for every x . ■

Example 3.6.3 $x \rightarrow d(x, y)$ is a continuous function from the metric space to the metric space of nonnegative real numbers.

This follows from Lemma 3.2.6. You can also define a metric on a Cartesian product of metric spaces.

Proposition 3.6.4 Let (X, d) be a metric space and consider $(X \times X, \rho)$ where

$$\rho((x, \tilde{x}), (y, \tilde{y})) \equiv d(x, y) + d(\tilde{x}, \tilde{y}).$$

Then this is also a metric space.

Proof: The only condition not obvious is the triangle inequality. However,

$$\begin{aligned} \rho((x, \tilde{x}), (y, \tilde{y})) + \rho((y, \tilde{y}), (z, \tilde{z})) &\equiv d(x, y) + d(\tilde{x}, \tilde{y}) + d(y, z) + d(\tilde{y}, \tilde{z}) \\ &\geq d(x, z) + d(\tilde{x}, \tilde{z}) = \rho((x, \tilde{x}), (z, \tilde{z})) \blacksquare \end{aligned}$$

Definition 3.6.5 If you have two metric spaces (X, d) and (Y, ρ) , a function $f : X \rightarrow Y$ is called a homeomorphism if and only if it is continuous, one to one, onto, and its inverse is also continuous.

Here is a useful proposition.

Proposition 3.6.6 Let (X, d) be a metric space and let S be a nonempty subset of X . Define

$$\text{dist}(x, S) \equiv \inf \{d(x, s) : s \in S\}$$

Then $|\text{dist}(x, S) - \text{dist}(y, S)| \leq d(x, y)$ so $x \rightarrow \text{dist}(x, S)$ is continuous.

Proof: Say $\text{dist}(x, S) \geq \text{dist}(y, S)$. Then there is $s \in S$ such that $\text{dist}(y, S) + \varepsilon > d(y, s)$. Then

$$\begin{aligned} |\text{dist}(x, S) - \text{dist}(y, S)| &= \text{dist}(x, S) - \text{dist}(y, S) \leq d(x, s) - (d(y, s) - \varepsilon) \\ &\leq d(x, y) + d(y, s) - (d(y, s) - \varepsilon) = d(x, y) + \varepsilon \end{aligned}$$

Since $\varepsilon > 0$ is arbitrary, this shows the claimed result. If $\text{dist}(x, S) \leq \text{dist}(y, S)$, repeat switching roles of x and y . ■

3.7 Continuity and Compactness

How does compactness relate to continuity? It turns out that the continuous image of a compact set is always compact. This is an easy consequence of the above major theorem.

Theorem 3.7.1 Let $f : X \rightarrow Y$ where (X, d) and (Y, ρ) are metric spaces and f is continuous on X . Then if $K \subseteq X$ is compact, it follows that $f(K)$ is compact in (Y, ρ) .

Proof: Let \mathcal{C} be an open cover of $f(K)$. Denote by $f^{-1}(\mathcal{C})$ the sets of the form $\{f^{-1}(U) : U \in \mathcal{C}\}$. Then $f^{-1}(\mathcal{C})$ is an open cover of K . It follows there are finitely many sets of the form $\{f^{-1}(U_1), \dots, f^{-1}(U_n)\}$ which covers K . It follows that $\{U_1, \dots, U_n\}$ is an open cover for $f(K)$. ■

The following is the important extreme values theorem for a real valued function defined on a compact set.

Theorem 3.7.2 *Let K be a compact metric space and suppose $f : K \rightarrow \mathbb{R}$ is a continuous function. That is, \mathbb{R} is the metric space where the metric is given by $d(x, y) = |x - y|$. Then f achieves its maximum and minimum values on K .*

Proof: Let $\lambda \equiv \sup \{f(x) : x \in K\}$. Then from the definition of sup, you have the existence of a sequence $\{x_n\} \subseteq K$ such that $\lim_{n \rightarrow \infty} f(x_n) = \lambda$. There is a subsequence still called $\{x_n\}$ which converges to some $x \in K$. From continuity, $\lambda = \lim_{n \rightarrow \infty} f(x_n) = f(x)$ and so f achieves its maximum value at x . Similar reasoning shows that it achieves its minimum value on K . ■

Definition 3.7.3 *Let $f : (X, d) \rightarrow (Y, \rho)$ be a function. Then it is said to be uniformly continuous on X if for every $\varepsilon > 0$ there exists a $\delta > 0$ such that whenever x, \hat{x} are two points of X with $d(x, \hat{x}) < \delta$, it follows that $\rho(f(x), f(\hat{x})) < \varepsilon$.*

Note the difference between this and continuity. With continuity, the δ could depend on x but here it works for any pair of points in X .

There is a remarkable result concerning compactness and uniform continuity.

Theorem 3.7.4 *Let $f : (X, d) \rightarrow (Y, \rho)$ be a continuous function and let K be a compact subset of X . Then the restriction of f to K is uniformly continuous.*

Proof: First of all, K is a metric space and f restricted to K is continuous. Now suppose it fails to be uniformly continuous. Then there exists $\varepsilon > 0$ and pairs of points x_n, \hat{x}_n such that $d(x_n, \hat{x}_n) < 1/n$ but $\rho(f(x_n), f(\hat{x}_n)) \geq \varepsilon$. Since K is compact, it is sequentially compact and so there exists a subsequence, still denoted as $\{x_n\}$ such that $x_n \rightarrow x \in K$. Then also $\hat{x}_n \rightarrow x$ also and so by Lemma 3.2.6, $\rho(f(x), f(x)) = \lim_{n \rightarrow \infty} \rho(f(x_n), f(\hat{x}_n)) \geq \varepsilon$ which is a contradiction. ■

3.8 Lipschitz Continuity and Contraction Maps

The following may be of more interest in the case of normed vector spaces, but there is no harm in stating it in this more general setting. You should verify that the functions described in the following definition are all continuous.

Definition 3.8.1 *Let $f : X \rightarrow Y$ where (X, d) and (Y, ρ) are metric spaces. Then f is said to be Lipschitz continuous if for every $x, \hat{x} \in X$, $\rho(f(x), f(\hat{x})) \leq rd(x, \hat{x})$. The function is called a contraction map if $r < 1$.*

The big theorem about contraction maps is the following.

Theorem 3.8.2 *Let $f : (X, d) \rightarrow (X, d)$ be a contraction map and let (X, d) be a complete metric space. Thus Cauchy sequences converge and also $d(f(x), f(\hat{x})) \leq rd(x, \hat{x})$ where $r < 1$. Then f has a unique fixed point. This is a point $x \in X$ such that $f(x) = x$. Also, if x_0 is any point of X , then*

$$d(x, x_0) \leq \frac{d(x_0, f(x_0))}{1 - r}$$

Also, for each n ,

$$d(f^n(x_0), x_0) \leq \frac{d(x_0, f(x_0))}{1 - r},$$

and $x = \lim_{n \rightarrow \infty} f^n(x_0)$.

Proof: Pick $x_0 \in X$ and consider the sequence of the iterates of the map f given by $x_0, f(x_0), f^2(x_0), \dots$. We argue that this is a Cauchy sequence. For $m < n$, it follows from the triangle inequality,

$$d(f^m(x_0), f^n(x_0)) \leq \sum_{k=m}^{n-1} d(f^{k+1}(x_0), f^k(x_0)) \leq \sum_{k=m}^{\infty} r^k d(f(x_0), x_0)$$

The reason for this last is as follows.

$$d(f^2(x_0), f(x_0)) \leq rd(f(x_0), x_0)$$

$$d(f^3(x_0), f^2(x_0)) \leq rd(f^2(x_0), f(x_0)) \leq r^2 d(f(x_0), x_0)$$

and so forth. Therefore, by the triangle inequality,

$$\begin{aligned} d(f^m(x_0), f^n(x_0)) &\leq \sum_{k=m}^{n-1} d(f^{k+1}(x_0), f^k(x_0)) \\ &\leq \sum_{k=m}^{\infty} r^k d(f(x_0), x_0) \leq d(f(x_0), x_0) \frac{r^m}{1-r} \end{aligned} \quad (3.1)$$

which shows that this is indeed a Cauchy sequence. Therefore, there exists x such that $\lim_{n \rightarrow \infty} f^n(x_0) = x$. By continuity, $f(x) = f(\lim_{n \rightarrow \infty} f^n(x_0)) = \lim_{n \rightarrow \infty} f^{n+1}(x_0) = x$.

Also note that, letting $m = 0$ in 3.1, this estimate yields

$$d(x_0, f^n(x_0)) \leq \frac{d(x_0, f(x_0))}{1-r}$$

Now $d(x_0, x) \leq d(x_0, f^n(x_0)) + d(f^n(x_0), x)$ and so

$$d(x_0, x) - d(f^n(x_0), x) \leq \frac{d(x_0, f(x_0))}{1-r}$$

Letting $n \rightarrow \infty$, it follows that $d(x_0, x) \leq \frac{d(x_0, f(x_0))}{1-r}$ because $\lim_{n \rightarrow \infty} d(f^n(x_0), x) = d(x, x) = 0$ by Lemma 3.2.6.

It only remains to verify that there is only one fixed point. Suppose then that x, x' are two. Then

$$d(x, x') = d(f(x), f(x')) \leq rd(x', x)$$

and so $d(x, x') = 0$ because $r < 1$. ■

The above is the usual formulation of this important theorem, but we actually proved a better result.

Corollary 3.8.3 *Let B be a closed subset of the complete metric space (X, d) and let $f : B \rightarrow X$ be a contraction map*

$$d(f(x), f(\hat{x})) \leq rd(x, \hat{x}), \quad r < 1.$$

*Also suppose **there exists** $x_0 \in B$ such that the sequence of iterates $\{f^n(x_0)\}_{n=1}^{\infty}$ remains in B . Then f has a unique fixed point in B which is the limit of the sequence of iterates. This is a point $x \in B$ such that $f(x) = x$. In the case that $B = B(x_0, \delta)$, the sequence of iterates satisfies the inequality*

$$d(f^n(x_0), x_0) \leq \frac{d(x_0, f(x_0))}{1-r}$$

and so it will remain in B if $\frac{d(x_0, f(x_0))}{1-r} < \delta$.

Proof: By assumption, the sequence of iterates stays in B . Then, as in the proof of the preceding theorem, for $m < n$, it follows from the triangle inequality,

$$\begin{aligned} d(f^m(x_0), f^n(x_0)) &\leq \sum_{k=m}^{n-1} d(f^{k+1}(x_0), f^k(x_0)) \\ &\leq \sum_{k=m}^{\infty} r^k d(f(x_0), x_0) = \frac{r^m}{1-r} d(f(x_0), x_0) \end{aligned}$$

Hence the sequence of iterates is Cauchy and must converge to a point x in X . However, B is closed and so it must be the case that $x \in B$. Then as before,

$$x = \lim_{n \rightarrow \infty} f^n(x_0) = \lim_{n \rightarrow \infty} f^{n+1}(x_0) = f\left(\lim_{n \rightarrow \infty} f^n(x_0)\right) = f(x)$$

As to the sequence of iterates remaining in B where B is a ball as described, the inequality above in the case where $m = 0$ yields $d(x_0, f^n(x_0)) \leq \frac{1}{1-r} d(f(x_0), x_0)$ and so, if the right side is less than δ , then the iterates remain in B . As to the fixed point being unique, it is as before. If x, x' are both fixed points in B , then $d(x, x') = d(f(x), f(x')) \leq rd(x, x')$ and so $x = x'$. ■

The contraction mapping theorem has an extremely useful generalization. In order to get a unique fixed point, it suffices to have some power of f a contraction map.

Theorem 3.8.4 *Let $f : (X, d) \rightarrow (X, d)$ have the property that for some $n \in \mathbb{N}$, f^n is a contraction map and let (X, d) be a complete metric space. Then there is a unique fixed point for f . As in the earlier theorem the sequence of iterates $\{f^n(x_0)\}_{n=1}^{\infty}$ also converges to the fixed point.*

Proof: From Theorem 3.8.2 there is a unique fixed point for f^n . Thus $f^n(x) = x$. Then

$$f^n(f(x)) = f^{n+1}(x) = f(x)$$

By uniqueness, $f(x) = x$.

Now consider the sequence of iterates. Suppose it fails to converge to x . Then there is $\varepsilon > 0$ and a subsequence n_k such that $d(f^{n_k}(x_0), x) \geq \varepsilon$. Now $n_k = p_k n + r_k$ where r_k is one of the numbers $\{0, 1, 2, \dots, n-1\}$. It follows that there exists one of these numbers which is repeated infinitely often. Call it r and let the further subsequence continue to be denoted as n_k . Thus $d(f^{p_k n + r}(x_0), x) \geq \varepsilon$. In other words,

$$d(f^{p_k n}(f^r(x_0)), x) \geq \varepsilon$$

However, from Theorem 3.8.2, as $k \rightarrow \infty$, $f^{p_k n}(f^r(x_0)) \rightarrow x$ which contradicts the above inequality. Hence the sequence of iterates converges to x , as it did for f a contraction map. ■

3.9 Convergence of Functions

Next is to consider the meaning of convergence of sequences of functions. There are two main ways of convergence of interest here, pointwise and uniform convergence.

Definition 3.9.1 *Let $f_n : X \rightarrow Y$ where $(X, d), (Y, \rho)$ are two metric spaces. Then $\{f_n\}$ is said to converge pointwise to a function $f : X \rightarrow Y$ if for every $x \in X$, $\lim_{n \rightarrow \infty} f_n(x) = f(x)$. $\{f_n\}$ is said to converge uniformly if for all $\varepsilon > 0$, there exists N such that if $n \geq N$, then $\sup_{x \in X} \rho(f_n(x), f(x)) < \varepsilon$*

Here is a well known example illustrating the difference between pointwise and uniform convergence.

Example 3.9.2 Let $f_n(x) = x^n$ on the metric space $[0, 1]$. Then this function converges pointwise to

$$f(x) = \begin{cases} 0 & \text{on } [0, 1) \\ 1 & \text{at } 1 \end{cases}$$

but it does not converge uniformly on this interval to f .

Note how the target function f in the above example is not continuous even though each function in the sequence is. The nice thing about uniform convergence is that it takes continuity of the functions in the sequence and imparts it to the target function. It does this for both continuity at a single point and uniform continuity. Thus uniform convergence is a very superior thing.

Theorem 3.9.3 Let $f_n : X \rightarrow Y$ where (X, d) , (Y, ρ) are two metric spaces and suppose each f_n is continuous at $x \in X$ and also that f_n converges uniformly to f on X . Then f is also continuous at x . In addition to this, if each f_n is uniformly continuous on X , then the same is true for f .

Proof: Let $\varepsilon > 0$ be given. Then

$$\rho(f(x), f(\hat{x})) \leq \rho(f(x), f_n(x)) + \rho(f_n(x), f_n(\hat{x})) + \rho(f_n(\hat{x}), f(\hat{x}))$$

By uniform convergence, there exists N such that both $\rho(f(x), f_n(x))$, $\rho(f_n(\hat{x}), f(\hat{x}))$ are less than $\varepsilon/3$ provided $n \geq N$. Thus picking such an n

$$\rho(f(x), f(\hat{x})) \leq \frac{2\varepsilon}{3} + \rho(f_n(x), f_n(\hat{x}))$$

From the continuity of f_n , there exists a positive number $\delta > 0$ such that if $d(x, \hat{x}) < \delta$, then $\rho(f_n(x), f_n(\hat{x})) < \varepsilon/3$. Hence, if $d(x, \hat{x}) < \delta$, then

$$\rho(f(x), f(\hat{x})) \leq \frac{2\varepsilon}{3} + \rho(f_n(x), f_n(\hat{x})) < \frac{2\varepsilon}{3} + \frac{\varepsilon}{3} = \varepsilon$$

Hence, f is continuous at x .

Next consider uniform continuity. It follows from the uniform convergence that if x, \hat{x} are any two points of X , then if $n \geq N$, then, picking such an n , $\rho(f(x), f(\hat{x})) \leq \frac{2\varepsilon}{3} + \rho(f_n(x), f_n(\hat{x}))$. By uniform continuity of f_n there exists δ such that if $d(x, \hat{x}) < \delta$, then the term on the right in the above is less than $\varepsilon/3$. Hence if $d(x, \hat{x}) < \delta$, then $\rho(f(x), f(\hat{x})) < \varepsilon$ and so f is uniformly continuous as claimed. ■

3.10 Compactness in $C(X, Y)$ Ascoli Arzela Theorem

This will use the characterization of compact metric spaces to give a proof of a general version of the Arzella Ascoli theorem. See Naylor and Sell [43] which is where I saw this general formulation.

Definition 3.10.1 Let (X, d_X) be a compact metric space. Let (Y, d_Y) be another complete metric space. Then $C(X, Y)$ will denote the continuous functions which map X to Y . Then ρ is a metric on $C(X, Y)$ defined by $\rho(f, g) \equiv \sup_{x \in X} d_Y(f(x), g(x))$.

Theorem 3.10.2 $(C(X, Y), \rho)$ is a complete metric space where (X, d_X) is a compact metric space

Proof: It is first necessary to show that ρ is well defined. In this argument, I will just write d rather than d_X or d_Y . To show this, note that from Lemma 3.2.6, if $x_n \rightarrow x$, and $y_n \rightarrow y$, then $d(x_n, y_n) \rightarrow d(x, y)$. Therefore, if f, g are continuous, and $x_n \rightarrow x$ so $f(x_n) \rightarrow f(x)$ and $g(x_n) \rightarrow g(x)$, $d(f(x_n), g(x_n)) \rightarrow d(f(x), g(x))$ and so, $\rho(f, g)$ is just the maximum of a continuous function defined on a compact set. By Theorem 3.7.2, the extreme values theorem, this maximum exists.

Clearly $\rho(f, g) = \rho(g, f)$ and

$$\begin{aligned} \rho(f, g) + \rho(g, h) &= \sup_{x \in X} d(f(x), g(x)) + \sup_{x \in X} d(g(x), h(x)) \\ &\geq \sup_{x \in X} (d(f(x), g(x)) + d(g(x), h(x))) \\ &\geq \sup_{x \in X} d(f(x), h(x)) = \rho(f, h) \end{aligned}$$

so the triangle inequality holds.

It remains to check completeness. Let $\{f_n\}$ be a Cauchy sequence. Then from the definition, $\{f_n(x)\}$ is a Cauchy sequence in Y and so it converges to something called $f(x)$. By Theorem 3.9.3, f is continuous. It remains to show that $\rho(f_n, f) \rightarrow 0$. Let $x \in X$. Then from what was just noted,

$$d(f_n(x), f(x)) = \lim_{m \rightarrow \infty} d(f_n(x), f_m(x)) \leq \limsup_{m \rightarrow \infty} \rho(f_n, f_m)$$

since $\{f_n\}$ is given to be a Cauchy sequence, there exists N such that if $m, n > N$, then $\rho(f_n, f_m) < \varepsilon$. Therefore, if $n > N$, $d(f_n(x), f(x)) \leq \limsup_{m \rightarrow \infty} \rho(f_n, f_m) \leq \varepsilon$. Since x is arbitrary, it follows that $\rho(f_n, f) \leq \varepsilon$, if $n \geq N$. ■

Here is a useful lemma.

Lemma 3.10.3 Let S be a totally bounded subset of (X, d) a metric space. Then \bar{S} is also totally bounded.

Proof: Suppose not. Then there exists a sequence $\{p_n\} \subseteq \bar{S}$ such that

$$d(p_m, p_n) \geq \varepsilon$$

for all $m \neq n$. Now let $q_n \in B(p_n, \frac{\varepsilon}{8}) \cap S$. Then it follows that

$$\frac{\varepsilon}{8} + d(q_n, q_m) + \frac{\varepsilon}{8} \geq d(p_n, q_n) + d(q_n, q_m) + d(q_m, p_m) \geq d(p_n, q_m) \geq \varepsilon$$

and so $d(q_n, q_m) > \frac{\varepsilon}{2}$. This contradicts total boundedness of S . ■

Next, here is an important definition.

Definition 3.10.4 Let $\mathcal{A} \subseteq C(X, Y)$ where (X, d_X) and (Y, d_Y) are metric spaces. Thus \mathcal{A} is a set of continuous functions mapping X to Y . Then \mathcal{A} is said to be equicontinuous if for every $\varepsilon > 0$ there exists a $\delta > 0$ such that if $d_X(x_1, x_2) < \delta$ then for all $f \in \mathcal{A}$, $d_Y(f(x_1), f(x_2)) < \varepsilon$. (This is uniform continuity which is uniform in \mathcal{A} .) \mathcal{A} is said to be pointwise compact if $\{f(x) : f \in \mathcal{A}\}$ has compact closure in Y .

Here is the Ascoli Arzela theorem.

Theorem 3.10.5 *Let (X, d_X) be a compact metric space and let (Y, d_Y) be a complete metric space. Thus $(C(X, Y), \rho)$ is a complete metric space. Let $\mathcal{A} \subseteq C(X, Y)$ be pointwise compact and equicontinuous. Then $\overline{\mathcal{A}}$ is compact. Here the closure is taken in $(C(X, Y), \rho)$.*

Proof: The more useful direction is that the two conditions imply compactness of $\overline{\mathcal{A}}$. I prove this first. Since $\overline{\mathcal{A}}$ is a closed subset of a complete space, it follows from Theorem 3.5.8, that $\overline{\mathcal{A}}$ will be compact if it is totally bounded. In showing this, it follows from Lemma 3.10.3 that it suffices to verify that \mathcal{A} is totally bounded. Suppose this is not so. Then there exists $\varepsilon > 0$ and a sequence of points of \mathcal{A} , $\{f_n\}$ such that $\rho(f_n, f_m) \geq \varepsilon$ whenever $n \neq m$.

By equicontinuity, there exists $\delta > 0$ such that if $d(x, y) < \delta$, then $d_Y(f(x), f(y)) < \frac{\varepsilon}{8}$ for all $f \in \mathcal{A}$. Let $\{x_i\}_{i=1}^p$ be a δ net for X . Since there are only finitely many x_i , it follows from pointwise compactness that there exists a subsequence, still denoted by $\{f_n\}$ which converges at each x_i . Now let $x \in X$ be arbitrary. There exists N such that for each x_i in that δ net,

$$d_Y(f_n(x_i), f_m(x_i)) < \varepsilon/8 \text{ whenever } n, m \geq N$$

Then for $m, n \geq N$,

$$\begin{aligned} & d_Y(f_n(x), f_m(x)) \\ & \leq d_Y(f_n(x), f_n(x_i)) + d_Y(f_n(x_i), f_m(x_i)) + d_Y(f_m(x_i), f_m(x)) \\ & < d_Y(f_n(x), f_n(x_i)) + \varepsilon/8 + d_Y(f_m(x_i), f_m(x)) \end{aligned}$$

Pick x_i such that $d(x, x_i) < \delta$. $\{x_i\}_{i=1}^p$ is a δ net and so this is surely possible. Then by equicontinuity, the two ends are each less than $\varepsilon/8$ and so for $m, n \geq N$,

$$d_Y(f_n(x), f_m(x)) \leq \frac{3\varepsilon}{8}$$

Since x is arbitrary, it follows that $\rho(f_n, f_m) \leq 3\varepsilon/8 < \varepsilon$ which is a contradiction. It follows that \mathcal{A} and hence $\overline{\mathcal{A}}$ is totally bounded. This proves the more important direction.

Next suppose $\overline{\mathcal{A}}$ is compact. Why must \mathcal{A} be pointwise compact and equicontinuous? If it fails to be pointwise compact, then there exists $x \in X$ such that $\{f(x) : f \in \mathcal{A}\}$ is not contained in a compact set of Y . Thus there exists $\varepsilon > 0$ and a sequence of functions in \mathcal{A} $\{f_n\}$ such that $d(f_n(x), f_m(x)) \geq \varepsilon$. But this implies $\rho(f_m, f_n) \geq \varepsilon$ and so $\overline{\mathcal{A}}$ fails to be totally bounded, a contradiction. Thus \mathcal{A} must be pointwise compact. Now why must it be equicontinuous? If it is not, then for each $n \in \mathbb{N}$ there exists $\varepsilon > 0$ and $x_n, y_n \in X$ such that $d(x_n, y_n) < 1/n$ but for some $f_n \in \mathcal{A}$, $d(f_n(x_n), f_n(y_n)) \geq \varepsilon$. However, by compactness, there exists a subsequence $\{f_{n_k}\}$ such that $\lim_{k \rightarrow \infty} \rho(f_{n_k}, f) = 0$ and also that $x_{n_k}, y_{n_k} \rightarrow x \in X$. Hence

$$\begin{aligned} \varepsilon & \leq d(f_{n_k}(x_{n_k}), f_{n_k}(y_{n_k})) \leq d(f_{n_k}(x_{n_k}), f(x_{n_k})) \\ & \quad + d(f(x_{n_k}), f(y_{n_k})) + d(f(y_{n_k}), f_{n_k}(y_{n_k})) \\ & \leq \rho(f_{n_k}, f) + d(f(x_{n_k}), f(y_{n_k})) + \rho(f, f_{n_k}) \end{aligned}$$

and now this is a contradiction because each term on the right converges to 0. The middle term converges to 0 because $f(x_{n_k}), f(y_{n_k}) \rightarrow f(x)$. See Lemma 3.2.6. ■

3.11 Connected Sets

Stated informally, connected sets are those which are in one piece. In order to define what is meant by this, I will first consider what it means for a set to **not** be in one piece. This is called **separated**. Connected sets are defined in terms of **not** being separated. This is why theorems about connected sets sometimes seem a little tricky.

Definition 3.11.1 A set, S in a metric space, is separated if there exist sets A, B such that

$$S = A \cup B, A, B \neq \emptyset, \text{ and } \bar{A} \cap B = \bar{B} \cap A = \emptyset.$$

In this case, the sets A and B are said to separate S . A set is connected if it is not separated. Remember \bar{A} denotes the closure of the set A .

Note that the concept of connected sets is defined in terms of what it is not. This makes it somewhat difficult to understand. One of the most important theorems about connected sets is the following.

Theorem 3.11.2 Suppose \mathcal{U} is a set of connected sets and that there exists a point p which is in all of these connected sets. Then $K \equiv \cup \mathcal{U}$ is connected.

Proof: The argument is dependent on Lemma 3.3.2. Suppose

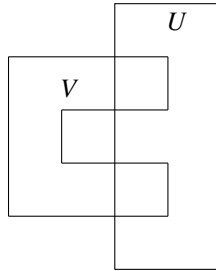
$$K = A \cup B$$

where $\bar{A} \cap B = \bar{B} \cap A = \emptyset, A \neq \emptyset, B \neq \emptyset$. Then p is in one of these sets. Say $p \in A$. Then if $U \in \mathcal{U}$, it must be the case that $U \subseteq A$ since if not, you would have

$$U = (A \cap U) \cup (B \cap U)$$

and the limit points of $A \cap U$ cannot be in B hence not in $B \cap U$ while the limit points of $B \cap U$ cannot be in A hence not in $A \cap U$. Thus $B = \emptyset$. It follows that K cannot be separated and so it is connected. ■

The intersection of connected sets is not necessarily connected as is shown by the following picture.



Theorem 3.11.3 Let $f : X \rightarrow Y$ be continuous where Y is a metric space and X is connected. Then $f(X)$ is also connected.

Proof: To do this you show $f(X)$ is not separated. Suppose to the contrary that $f(X) = A \cup B$ where A and B separate $f(X)$. Then consider the sets $f^{-1}(A)$ and $f^{-1}(B)$. If $z \in f^{-1}(B)$, then $f(z) \in B$ and so $f(z)$ is not a limit point of A . Therefore, there exists an

open set, U containing $f(z)$ such that $U \cap A = \emptyset$. But then, the continuity of f and Theorem 3.6.2 implies that $f^{-1}(U)$ is an open set containing z such that $f^{-1}(U) \cap f^{-1}(A) = \emptyset$. Therefore, $f^{-1}(B)$ contains no limit points of $f^{-1}(A)$. Similar reasoning implies $f^{-1}(A)$ contains no limit points of $f^{-1}(B)$. It follows that X is separated by $f^{-1}(A)$ and $f^{-1}(B)$, contradicting the assumption that X was connected. ■

An arbitrary set can be written as a union of maximal connected sets called connected components. This is the concept of the next definition.

Definition 3.11.4 Let S be a set and let $p \in S$. Denote by C_p the union of all connected subsets of S which contain p . This is called the connected component determined by p .

Theorem 3.11.5 Let C_p be a connected component of a set S in a metric space. Then C_p is a connected set and if $C_p \cap C_q \neq \emptyset$, then $C_p = C_q$.

Proof: Let \mathcal{C} denote the connected subsets of S which contain p . By Theorem 3.11.2, $\bigcup \mathcal{C} = C_p$ is connected. If $x \in C_p \cap C_q$, then from Theorem 3.11.2, $C_p \supseteq C_p \cup C_q$ and so $C_p \supseteq C_q$. The inclusion goes the other way by the same reason. ■

This shows the connected components of a set are equivalence classes and partition the set.

A set, I is an interval in \mathbb{R} if and only if whenever $x, y \in I$ then $[x, y] \subseteq I$. The following theorem is about the connected sets in \mathbb{R} .

Theorem 3.11.6 A set C in \mathbb{R} is connected if and only if C is an interval.

Proof: Let C be connected. If C consists of a single point, p , there is nothing to prove. The interval is just $[p, p]$. Suppose $p < q$ and $p, q \in C$. You need to show $(p, q) \subseteq C$. If

$$x \in (p, q) \setminus C$$

let $C \cap (-\infty, x) \equiv A$, and $C \cap (x, \infty) \equiv B$. Then $C = A \cup B$ and the sets A and B separate C contrary to the assumption that C is connected.

Conversely, let I be an interval. Suppose I is separated by A and B . Pick $x \in A$ and $y \in B$. Suppose without loss of generality that $x < y$. Now define the set,

$$S \equiv \{t \in [x, y] : [x, t] \subseteq A\}$$

and let l be the least upper bound of S . Then $l \in \bar{A}$ so $l \notin B$ which implies $l \in A$. But if $l \notin \bar{B}$, then for some $\delta > 0$,

$$(l, l + \delta) \cap B = \emptyset$$

contradicting the definition of l as an upper bound for S . Therefore, $l \in \bar{B}$ which implies $l \notin A$ after all, a contradiction. It follows I must be connected. ■

This yields a generalization of the intermediate value theorem from one variable calculus.

Corollary 3.11.7 Let E be a connected set in a metric space and suppose $f : E \rightarrow \mathbb{R}$ and that $y \in (f(e_1), f(e_2))$ where $e_i \in E$. Then there exists $e \in E$ such that $f(e) = y$.

Proof: From Theorem 3.11.3, $f(E)$ is a connected subset of \mathbb{R} . By Theorem 3.11.6 $f(E)$ must be an interval. In particular, it must contain y . This proves the corollary. ■

The following theorem is a very useful description of the open sets in \mathbb{R} .

Theorem 3.11.8 *Let U be an open set in \mathbb{R} . Then there exist countably many disjoint open sets $\{(a_i, b_i)\}_{i=1}^{\infty}$ such that $U = \cup_{i=1}^{\infty} (a_i, b_i)$.*

Proof: Let $p \in U$ and let $z \in C_p$, the connected component determined by p . Since U is open, there exists, $\delta > 0$ such that $(z - \delta, z + \delta) \subseteq U$. It follows from Theorem 3.11.2 that $(z - \delta, z + \delta) \subseteq C_p$. This shows C_p is open. By Theorem 3.11.6, this shows C_p is an open interval, (a, b) where $a, b \in [-\infty, \infty]$. There are therefore at most countably many of these connected components because each must contain a rational number and the rational numbers are countable. Denote by $\{(a_i, b_i)\}_{i=1}^{\infty}$ the set of these connected components. ■

Definition 3.11.9 *A set E in a metric space is arcwise connected if for any two points, $p, q \in E$, there exists a closed interval, $[a, b]$ and a continuous function, $\gamma : [a, b] \rightarrow E$ such that $\gamma(a) = p$ and $\gamma(b) = q$.*

An example of an arcwise connected metric space would be any subset of \mathbb{R}^n which is the continuous image of an interval. Arcwise connected is not the same as connected. A well known example is the following.

$$\left\{ \left(x, \sin \frac{1}{x} \right) : x \in (0, 1] \right\} \cup \{(0, y) : y \in [-1, 1]\} \quad (3.2)$$

You can verify that this set of points in the normed vector space \mathbb{R}^2 is not arcwise connected but is connected.

Lemma 3.11.10 *In \mathbb{R}^p , $B(z, r)$ is arcwise connected.*

Proof: This is easy from the convexity of the set. If $x, y \in B(z, r)$, then let $\gamma(t) = x + t(y - x)$ for $t \in [0, 1]$.

$$\begin{aligned} \|x + t(y - x) - z\| &= \|(1 - t)(x - z) + t(y - z)\| \\ &\leq (1 - t)\|x - z\| + t\|y - z\| \\ &< (1 - t)r + tr = r \end{aligned}$$

showing $\gamma(t)$ stays in $B(z, r)$. ■

Proposition 3.11.11 *If $X \neq \emptyset$ is arcwise connected, then it is connected.*

Proof: Let $p \in X$. Then by assumption, for any $x \in X$, there is an arc joining p and x . This arc is connected because it is the continuous image of an interval which is connected. Since x is arbitrary, every x is in a connected subset of X which contains p . Hence $C_p = X$ and so X is connected. ■

Theorem 3.11.12 *Let U be an open subset of \mathbb{R}^p . Then U is arcwise connected if and only if U is connected. Also the connected components of an open set are open sets.*

Proof: By Proposition 3.11.11 it is only necessary to verify that if U is connected and open, then U is arcwise connected. Pick $p \in U$. Say $x \in U$ satisfies \mathcal{P} if there exists a continuous function, $\gamma : [a, b] \rightarrow U$ such that $\gamma(a) = p$ and $\gamma(b) = x$.

$$A \equiv \{x \in U \text{ such that } x \text{ satisfies } \mathcal{P}\}$$

If $x \in A$, then Lemma 3.11.10 implies $B(x, r) \subseteq U$ is arcwise connected for small enough r . Thus letting $y \in B(x, r)$, there exist intervals, $[a, b]$ and $[c, d]$ and continuous functions having values in U , γ, η such that $\gamma(a) = p, \gamma(b) = x, \eta(c) = x$, and $\eta(d) = y$. Then let $\gamma_1 : [a, b + d - c] \rightarrow U$ be defined as

$$\gamma_1(t) \equiv \begin{cases} \gamma(t) & \text{if } t \in [a, b] \\ \eta(t + c - b) & \text{if } t \in [b, b + d - c] \end{cases}$$

Then it is clear that γ_1 is a continuous function mapping p to y and showing that $B(x, r) \subseteq A$. Therefore, A is open. $A \neq \emptyset$ because since U is open there is an open set, $B(p, \delta)$ containing p which is contained in U and is arcwise connected.

Now consider $B \equiv U \setminus A$. I claim this is also open. If B is not open, there exists a point $z \in B$ such that every open set containing z is not contained in B . Therefore, letting $B(z, \delta)$ be such that $z \in B(z, \delta) \subseteq U$, there exist points of A contained in $B(z, \delta)$. But then, a repeat of the above argument shows $z \in A$ also. Hence B is open and so if $B \neq \emptyset$, then $U = B \cup A$ and so U is separated by the two sets B and A contradicting the assumption that U is connected. Note that, since B is open, it contains no limit points of A and since A is open, it contains no limit points of B .

It remains to verify the connected components are open. Let $z \in C_p$ where C_p is the connected component determined by p . Then picking $B(z, \delta) \subseteq U$, $C_p \cup B(z, \delta)$ is connected and contained in U and so it must also be contained in C_p . Thus z is an interior point of C_p . ■

As an application, consider the following corollary.

Corollary 3.11.13 *Let $f : \Omega \rightarrow \mathbb{Z}$ be continuous where Ω is a connected nonempty open set of a metric space. Then f must be a constant.*

Proof: Suppose not. Then it achieves two different values, k and $l \neq k$. Then $\Omega = f^{-1}(l) \cup f^{-1}(\{m \in \mathbb{Z} : m \neq l\})$ and these are disjoint nonempty open sets which separate Ω . To see they are open, note

$$f^{-1}(\{m \in \mathbb{Z} : m \neq l\}) = f^{-1}\left(\bigcup_{m \neq l} \left(m - \frac{1}{6}, m + \frac{1}{6}\right)\right)$$

which is the inverse image of an open set while $f^{-1}(l) = f^{-1}\left((l - \frac{1}{6}, l + \frac{1}{6})\right)$ also an open set. ■

3.12 Partitions of Unity in Metric Space

Lemma 3.12.1 *Let X be a metric space and let S be a nonempty subset of X .*

$$\text{dist}(x, S) \equiv \inf\{d(x, z) : z \in S\}$$

Then

$$|\text{dist}(x, S) - \text{dist}(y, S)| \leq d(x, y).$$

Proof: Say $\text{dist}(x, S) \geq \text{dist}(y, S)$. Then letting $\varepsilon > 0$ be given, there exists $z \in S$ such that $d(y, z) < \text{dist}(y, S) + \varepsilon$. Then

$$|\text{dist}(x, S) - \text{dist}(y, S)| = \text{dist}(x, S) - \text{dist}(y, S) \leq \text{dist}(x, S) - (d(y, z) - \varepsilon)$$

$$\leq d(x, z) - (d(y, z) - \varepsilon) \leq d(x, y) + d(y, z) - d(y, z) + \varepsilon = d(x, y) + \varepsilon$$

Since ε is arbitrary, $|\text{dist}(x, S) - \text{dist}(y, S)| \leq d(x, y)$. The situation is completely similar if $\text{dist}(x, S) < \text{dist}(y, S)$. ■

Then this shows that $x \rightarrow \text{dist}(x, S)$ is a continuous real valued function.

This is about partitions of unity in metric space. Assume here that closed balls are compact. For example, you might be considering \mathbb{R}^p with $d(x, y) \equiv |x - y|$.

Definition 3.12.2 Define $\text{spt}(f)$ (support of f) to be the closure of the set $\{x : f(x) \neq 0\}$. If V is an open set, $C_c(V)$ will be the set of continuous functions f , defined on Ω having $\text{spt}(f) \subseteq V$.

Definition 3.12.3 If K is a compact subset of an open set, V , then $K \prec \phi \prec V$ if $\phi \in C_c(V)$, $\phi(K) = \{1\}$, $\phi(\Omega) \subseteq [0, 1]$, where Ω denotes the whole metric space. Also for $\phi \in C_c(\Omega)$, $K \prec \phi$ if $\phi(\Omega) \subseteq [0, 1]$ and $\phi(K) = 1$. $\phi \prec V$ if $\phi(\Omega) \subseteq [0, 1]$ and $\text{spt}(\phi) \subseteq V$.

Lemma 3.12.4 Let (Ω, d) be a metric space in which closed balls are compact. Then if K is a compact subset of an open set V , then there exists ϕ such that $K \prec \phi \prec V$.

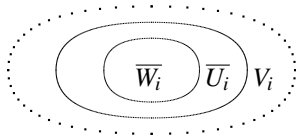
Proof: Since K is compact, the distance between K and V^C is positive, $\delta > 0$. Otherwise there would be $x_n \in K$ and $y_n \in V^C$ with $d(x_n, y_n) < 1/n$. Taking a subsequence, still denoted with n , we can assume $x_n \rightarrow x$ and $y_n \rightarrow x$ but this would imply x is in both K and V^C which is not possible. Now consider $\{B(x, \delta/2)\}$ for $x \in K$. This is an open cover and the closure of each ball is contained in V . Since K is compact, finitely many of these balls cover K . Denote their union as W . Then \bar{W} is compact because it is the finite union of the closed balls. Hence $K \subseteq W \subseteq \bar{W} \subseteq V$. Now consider

$$\phi(x) \equiv \frac{\text{dist}(x, W^C)}{\text{dist}(x, K) + \text{dist}(x, W^C)}$$

the denominator is never zero because x cannot be in both K and W^C . Thus ϕ is continuous by Lemma 3.12.1. also if $x \in K$, then $\phi(x) = 1$ and if $x \notin W$, then $\phi(x) = 0$. ■

Theorem 3.12.5 (Partition of unity) Let K be a compact subset of a metric space in which closed balls are compact and suppose $K \subseteq V = \bigcup_{i=1}^n V_i$, V_i open. Then there exist $\psi_i \prec V_i$ with $\sum_{i=1}^n \psi_i(x) = 1$ for all $x \in K$.

Proof: Let $K_1 = K \setminus \bigcup_{i=2}^n V_i$. Thus K_1 is compact and $K_1 \subseteq V_1$. Let $K_1 \subseteq W_1 \subseteq \bar{W}_1 \subseteq V_1$ with \bar{W}_1 compact. To obtain W_1 , use Lemma 3.12.4 to get f such that $K_1 \prec f \prec V_1$ and let $W_1 \equiv \{x : f(x) \neq 0\}$. Thus W_1, V_2, \dots, V_n covers K and $\bar{W}_1 \subseteq V_1$. Let $K_2 = K \setminus (\bigcup_{i=2}^n V_i \cup W_1)$. Then K_2 is compact and $K_2 \subseteq V_2$. Let $K_2 \subseteq W_2 \subseteq \bar{W}_2 \subseteq V_2$, \bar{W}_2 compact. Continue this way finally obtaining W_1, \dots, W_n , $K \subseteq W_1 \cup \dots \cup W_n$, and $\bar{W}_i \subseteq V_i$, \bar{W}_i compact. Now let $\bar{W}_i \subseteq U_i \subseteq \bar{U}_i \subseteq V_i$, \bar{U}_i compact.



By Lemma 3.12.4, let $\bar{U}_i \prec \phi_i \prec V_i$, $\cup_{i=1}^n \bar{W}_i \prec \gamma \prec \cup_{i=1}^n U_i$. Define

$$\psi_i(x) = \begin{cases} \gamma(x)\phi_i(x)/\sum_{j=1}^n \phi_j(x) & \text{if } \sum_{j=1}^n \phi_j(x) \neq 0, \\ 0 & \text{if } \sum_{j=1}^n \phi_j(x) = 0. \end{cases}$$

If x is such that $\sum_{j=1}^n \phi_j(x) = 0$, then $x \notin \cup_{i=1}^n \bar{U}_i$. Consequently $\gamma(y) = 0$ for all y near x and so $\psi_i(y) = 0$ for all y near x . Hence ψ_i is continuous at such x . If $\sum_{j=1}^n \phi_j(x) \neq 0$, this situation persists near x and so ψ_i is continuous at such points. Therefore ψ_i is continuous. If $x \in K$, then $\gamma(x) = 1$ and so $\sum_{j=1}^n \psi_j(x) = 1$. Clearly $0 \leq \psi_i(x) \leq 1$ and $\text{spt}(\psi_j) \subseteq V_j$. ■

3.13 Completion of Metric Spaces

Let (X, d) be a metric space $X \neq \emptyset$. Perhaps this is not a complete metric space. In other words, it may be that Cauchy Sequences do not converge. Of course if $x \in X$ and if $x_n = x$ for all n then $\{x_n\}$ is a Cauchy sequence and it converges to x .

Lemma 3.13.1 *Denote by x a Cauchy sequence x being short for $\{x_n\}_{n=1}^\infty$. Then if x, y are two Cauchy sequences, $\lim_{n \rightarrow \infty} d(x_n, y_n)$ exists.*

Proof: Let $\varepsilon > 0$ be given and let N be so large that whenever $n, m \geq N$, it follows that $d(x_n, x_m), d(y_n, y_m) < \varepsilon/2$. Then for such n, m

$$\begin{aligned} |d(x_n, y_n) - d(x_m, y_m)| &\leq |d(x_n, y_n) - d(x_n, y_m)| + |d(x_n, y_m) - d(x_m, y_m)| \\ &\leq d(y_n, y_m) + d(x_n, x_m) < \varepsilon \end{aligned}$$

by Lemma 3.12.1. Therefore, $\{d(x_n, y_n)\}_n$ is a Cauchy sequence in \mathbb{R} and so it converges. ■

Definition 3.13.2 *Let $x \sim y$ when $\lim_{n \rightarrow \infty} d(x_n, y_n) = 0$.*

Lemma 3.13.3 *\sim is an equivalence relation.*

Proof: Clearly $x \sim x$ and if $x \sim y$ then $y \sim x$. Suppose then that $x \sim y$ and $y \sim z$. Is $x \sim z$?

$$d(x_n, z_n) \leq d(x_n, y_n) + d(y_n, z_n)$$

and both of those terms on the right converge to 0. ■

Definition 3.13.4 *Denote by $[x]$ the equivalence class determined by the Cauchy sequence x . Let $d([x], [y]) \equiv \lim_{n \rightarrow \infty} d(x_n, y_n)$.*

Theorem 3.13.5 *Denote by \hat{X} the set of equivalence classes. Then d defined above is a metric, \hat{X} with this is a complete metric space, and X can be considered a dense subset of \hat{X} .*

Proof: That d just defined is a metric is obvious from the fact that the original metric d satisfies the triangle inequality. It is also clear that $d([x], [y]) \geq 0$ and that if $[x] = [y]$ if and only if $d([x], [y]) = 0$.

It remains to show that (\hat{X}, d) is complete. Let $\{[x]_n\}_n$ be a Cauchy sequence. From Theorem 3.2.2 it suffices to show the convergence of a subsequence. There is a subsequence, denoted as $\{[x^n]\}$ where x^n is a representative of $[x]_n$ such that $d([x^n], [x^{n+1}]) <$

4^{-n} . Thus there is an increasing sequence $\{k_n\}$ such that $d(x_k^n, x_l^{n+1}) < 2^{-n}$ if $k, l \geq k_n$ where k_n is increasing in n . Let $\mathbf{y} = \{x_{k_n}^n\}_{n=1}^\infty$. For $m \geq k_n$ and the triangle inequality,

$$\begin{aligned} d(x_m^n, y_m) &= d(x_m^n, x_{k_m}^m) \leq d(x_m^n, x_{k_n}^n) + d(x_{k_n}^n, x_{k_m}^m) \leq 2^{-n} + \sum_{j=n}^{m-1} d(x_{k_j}^j, x_{k_m}^{j+1}) \\ &< 2^{-n} + \sum_{j=n}^{m-1} 2^{-j} < 2^{-n} + 2^{-(n-1)} < 2^{-(n-2)} \end{aligned}$$

Then \mathbf{y} is a Cauchy sequence since it is a subsequence of one and also $d([x^n], [\mathbf{y}]) \rightarrow 0$.

To show that X is dense in \hat{X} , let $[x]$ be given. Then for m large enough, $d(x_k, x_m) < \varepsilon$ whenever $k \geq m$. It suffices to let \mathbf{y} be the constant Cauchy sequence always equal to x_m . ■

3.14 Exercises

1. Let $d(x, y) = |x - y|$ for $x, y \in \mathbb{R}$. Show that this is a metric on \mathbb{R} .
2. Now consider \mathbb{R}^n . Let $\|\mathbf{x}\|_\infty \equiv \max\{|x_i|, i = 1, \dots, n\}$. Define $d(\mathbf{x}, \mathbf{y}) \equiv \|\mathbf{x} - \mathbf{y}\|_\infty$. Show that this is a metric on \mathbb{R}^n . In the case of $n = 2$, describe the ball $B(\mathbf{0}, r)$. **Hint:** First show that $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$.
3. Let $C([0, T])$ denote the space of functions which are continuous on $[0, T]$. Define

$$\|f\| \equiv \|f\|_\infty \equiv \sup_{t \in [0, T]} |f(t)| = \max_{t \in [0, T]} |f(t)|$$

Verify the following. $\|f + g\| \leq \|f\| + \|g\|$. Then use to show that $d(f, g) \equiv \|f - g\|$ is a metric and that with this metric, $(C([0, T]), d)$ is a metric space.

4. Recall that $[a, b]$ is compact. Also, it is Lemma 3.5.9 above. Thus every open cover has a finite subcover of the set. Also recall that a sequence of numbers $\{x_n\}$ is a Cauchy sequence means that for every $\varepsilon > 0$ there exists N such that if $m, n > N$, then $|x_n - x_m| < \varepsilon$. First show that every Cauchy sequence is bounded. Next, using the compactness of closed intervals, show that every Cauchy sequence has a convergent subsequence. By Theorem 3.2.2, the original Cauchy sequence converges. Thus \mathbb{R} with the usual metric just described is complete because every Cauchy sequence converges.
5. Using the result of the above problem, show that $(\mathbb{R}^n, \|\cdot\|_\infty)$ is a complete metric space. That is, every Cauchy sequence converges. Here $d(\mathbf{x}, \mathbf{y}) \equiv \|\mathbf{x} - \mathbf{y}\|_\infty$.
6. Suppose you had (X_i, d_i) is a metric space. Now consider the product space $X \equiv \prod_{i=1}^n X_i$ with $d(\mathbf{x}, \mathbf{y}) = \max\{d(x_i, y_i), i = 1 \dots, n\}$. Would this be a metric space? If so, prove that this is the case.

Does triangle inequality hold? **Hint:** For each i ,

$$d_i(x_i, z_i) \leq d_i(x_i, y_i) + d_i(y_i, z_i) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$$

Now take max of the two ends.

7. In the above example, if each (X_i, d_i) is complete, explain why (X, d) is also complete.
8. Show that $C([0, T])$ is a complete metric space. That is, show that if $\{f_n\}$ is a Cauchy sequence, then there exists $f \in C([0, T])$ such that

$$\lim_{n \rightarrow \infty} d(f, f_n) = \lim_{n \rightarrow \infty} \|f - f_n\| = 0$$

This is just a special case of theorems discussed in the chapter.

9. Let X be a nonempty set of points. Say it has infinitely many points. Define $d(x, y) = 1$ if $x \neq y$ and $d(x, y) = 0$ if $x = y$. Show that this is a metric. Show that in (X, d) every point is open and closed. In fact, show that every set is open and every set is closed. Is this a complete metric space? Explain why. Describe the open balls.
10. Show that the union of any set of open sets is an open set. Show the intersection of any set of closed sets is closed. Let A be a nonempty subset of a metric space (X, d) . Then the closure of A , written as \bar{A} is defined to be the intersection of all closed sets which contain A . Show that $\bar{A} = A \cup A'$. That is, to find the closure, you just take the set and include all limit points of the set. It was proved in the chapter, but go over it yourself.
11. Let A' denote the set of limit points of A , a nonempty subset of a metric space (X, d) . Show that A' is closed.
12. A theorem was proved which gave three equivalent descriptions of compactness of a metric space. One of them said the following: A metric space is compact if and only if it is complete and totally bounded. Suppose (X, d) is a complete metric space and $K \subseteq X$. Then (K, d) is also clearly a metric space having the same metric as X . Show that (K, d) is compact if and only if it is **closed** and totally bounded. Note the similarity with the Heine Borel theorem on \mathbb{R} . Show that on \mathbb{R} , every bounded set is also totally bounded. Thus the earlier Heine Borel theorem for \mathbb{R} is obtained.
13. Suppose (X_i, d_i) is a compact metric space. Then the Cartesian product is also a metric space. That is $(\prod_{i=1}^n X_i, d)$ is a metric space where $d(\mathbf{x}, \mathbf{y}) \equiv \max \{d_i(x_i, y_i)\}$. Show that $(\prod_{i=1}^n X_i, d)$ is compact. Recall the Heine Borel theorem for \mathbb{R} . Explain why $\prod_{i=1}^n [a_i, b_i]$ is compact in \mathbb{R}^n with the distance given by

$$d(\mathbf{x}, \mathbf{y}) = \max \{|x_i - y_i|\}$$

Hint: It suffices to show that $(\prod_{i=1}^n X_i, d)$ is sequentially compact. Let $\{\mathbf{x}^m\}_{m=1}^\infty$ be a sequence. Then $\{x_1^m\}_{m=1}^\infty$ is a sequence in X_1 . Therefore, it has a subsequence $\{x_1^{k_1}\}_{k_1=1}^\infty$ which converges to a point $x_1 \in X_1$. Now consider $\{x_2^{k_1}\}_{k_1=1}^\infty$ the second components. It has a subsequence denoted as k_2 such that $\{x_2^{k_2}\}_{k_2=1}^\infty$ converges to a point x_2 in X_2 . Explain why $\lim_{k_2 \rightarrow \infty} x_1^{k_2} = x_1$. Continue doing this n times. Explain why $\lim_{k_n \rightarrow \infty} x_l^{k_n} = x_l \in X_l$ for each l . Then explain why this is the same as saying $\lim_{k_n \rightarrow \infty} \mathbf{x}^{k_n} = \mathbf{x}$ in $(\prod_{i=1}^n X_i, d)$.

14. If you have a metric space (X, d) and a compact subset of (X, d) K , suppose that L is a closed subset of K . Explain why L must also be compact. **Hint:** Use the definition of compactness. Explain why every closed and bounded set in \mathbb{R}^n is compact. Here the distance is given by $d(\mathbf{x}, \mathbf{y}) \equiv \max_{1 \leq i \leq n} \{|x_i - y_i|\}$.
15. Show that compactness is a topological property. If $(X, d), (Y, \rho)$ are both metric spaces and $f : X \rightarrow Y$ has the property that f is one to one, onto, and continuous, and also f^{-1} is one to one onto and continuous, then the two metric spaces are compact or not compact together. That is one is compact if and only if the other is.
16. Consider \mathbb{R} the real numbers. Define a distance in the following way. $\rho(x, y) \equiv |\arctan(x) - \arctan(y)|$. Show this is a good enough distance and that the open sets which come from this distance are the same as the open sets which come from the usual distance $d(x, y) = |x - y|$. Explain why this yields that the identity mapping $f(x) = x$ is continuous with continuous inverse as a map from (\mathbb{R}, d) to (\mathbb{R}, ρ) . To do this, you show that an open ball taken with respect to one of these is also open with respect to the other. However, (\mathbb{R}, ρ) is not a complete metric space while (\mathbb{R}, d) is. Thus, unlike compactness. Completeness is not a topological property. **Hint:** To show the lack of completeness of (\mathbb{R}, ρ) , consider $x_n = n$. Show it is a Cauchy sequence with respect to ρ .
17. If K is a compact subset of (X, d) and $y \notin K$, show that there always exists $x \in K$ such that $d(x, y) = \text{dist}(y, K)$. Give an example in \mathbb{R} to show that this might not be so if K is not compact.
18. If S is a nonempty set, the diameter of S denoted as $\text{diam}(S)$ is defined as follows. $\text{diam}(S) \equiv \sup \{d(x, y) : x, y \in S\}$. Suppose (X, d) is a complete metric space and you have a nested sequence of closed sets whose diameters converge to 0. That is, each A_n is closed, $\cdots A_n \supseteq A_{n+1} \cdots$ and $\lim_{n \rightarrow \infty} \text{diam}(A_n) = 0$. Show that there is exactly one point p contained in the intersection of all these sets A_n . Give an example which shows that if the condition on the diameters does not hold, then maybe there is no point in the intersection of these sets.
19. Two metric spaces $(X, d), (Y, \rho)$ are homeomorphic if there exists a continuous function $f : X \rightarrow Y$ which is one to one onto, and whose inverse is also continuous one to one and onto. Show that the interval $[0, 1]$ is not homeomorphic to the unit circle. **Hint:** Recall that the continuous image of a connected set is connected, Theorem 3.11.3. However, if you remove a point from $[0, 1]$ it is no longer connected but removing a single point from the circle results in a connected set.
20. Using the same methods in the above problem, show that the unit circle is not homeomorphic to the unit sphere $\{x^2 + y^2 + z^2 = 1\}$ and the unit circle is not homeomorphic to a figure eight.
21. The rational numbers \mathbb{Q} and the natural numbers \mathbb{N} have the property that there is a one to one and onto map from \mathbb{N} to \mathbb{Q} . This is a simple consequence of the Schroeder Bernstein theorem presented earlier. Both of these are also metric spaces with respect to the usual metric on \mathbb{R} . Are they homeomorphic? **Hint:** Suppose they were. Then in \mathbb{Q} consider $(1, 2)$, all the rationals between 1 and 2 excluding 1 and 2. This is not a closed set because 2 is a limit point of the set which is not in it. Now if you have f a homeomorphism, consider $f((1, 2))$. Is this set closed?

22. If you have an open set O in \mathbb{R} , show that O is the countable union of disjoint open intervals. **Hint:** Consider the connected components. Go over this for yourself. It is in the chapter.
23. Addition and multiplication on \mathbb{R} can be considered mappings from $\mathbb{R} \times \mathbb{R}$ to \mathbb{R} as follows. $+(x, y) \equiv x + y$, $\cdot(x, y) \equiv xy$. Here the metric on $\mathbb{R} \times \mathbb{R}$ can be taken as $d((x, y), (\hat{x}, \hat{y})) \equiv \max(|x - \hat{x}|, |y - \hat{y}|)$. Show these operations are continuous functions.
24. Suppose K is a compact subset of a metric space (X, d) and there is an open cover \mathcal{C} of K . Show that there exists a single positive $\delta > 0$ such that if $x \in K$, $B(x, \delta)$ is contained in some set of \mathcal{C} . This number is called a Lebesgue number. Do this directly from the definition of compactness in terms of open covers without using the equivalence of compactness and sequential compactness.
25. Show uniform continuity of a continuous function defined on a compact set where compactness only refers to open covers. Use the above problem on existence of the Lebesgue number.
26. Let $f : D \rightarrow \mathbb{R}$ be a function. This function is said to be lower semicontinuous¹ at $x \in D$ if for any sequence $\{x_n\} \subseteq D$ which converges to x it follows $f(x) \leq \liminf_{n \rightarrow \infty} f(x_n)$. Suppose D is sequentially compact and f is lower semicontinuous at every point of D . Show that then f achieves its minimum on D . Here D is some metric space. Let $f : D \rightarrow \mathbb{R}$ be a function. This function is said to be upper semicontinuous at $x \in D$ if for any sequence $\{x_n\} \subseteq D$ which converges to x it follows $f(x) \geq \limsup_{n \rightarrow \infty} f(x_n)$. Suppose D is sequentially compact and f is upper semicontinuous at every point of D . Show that then f achieves its maximum on D .
27. Show that a real valued function defined on a metric space D is continuous if and only if it is both upper and lower semicontinuous.
28. Give an example of a lower semicontinuous function defined on \mathbb{R} which is not continuous and an example of an upper semicontinuous function which is not continuous.
29. More generally, one considers functions which have values in $[-\infty, \infty]$. Then f is upper semicontinuous if, whenever $x_n \rightarrow x$, $f(x) \geq \limsup_{n \rightarrow \infty} f(x_n)$ and lower semicontinuous if whenever $x_n \rightarrow x$, $f(x) \leq \liminf_{n \rightarrow \infty} f(x_n)$. Suppose $\{f_\alpha : \alpha \in \Lambda\}$ is a collection of continuous real valued functions defined on a metric space. Let $F(x) \equiv \inf\{f_\alpha(x) : \alpha \in \Lambda\}$. Show F is an upper semicontinuous function. Next let $G(x) \equiv \sup\{f_\alpha(x) : \alpha \in \Lambda\}$. Show G is a lower semicontinuous function.
30. The result of this problem is due to Hausdorff. It says that if you have any lower semicontinuous real valued function defined on a metric space (X, d) , then it is the limit of an increasing sequence of continuous functions. Here is an outline. You complete the details.
- (a) First suppose $f(x) \geq 0$ for all x . Define $f_n(x) \equiv \inf_{z \in X} \{f(z) + nd(z, x)\}$. Then $f(x) \geq f_n(x)$ and $f_n(x)$ is increasing in n . Also each f_n is continuous because

¹The notion of lower semicontinuity is very important for functions which are defined on infinite dimensional sets.

$f_n(x) \leq f(z) + nd(z, y) + nd(y, x)$. Thus $f_n(x) \leq f_n(y) + nd(y, x)$. Why? It follows that $|f_n(x) - f_n(y)| \leq nd(y, x)$. Why?

- (b) Let $h(x) = \lim_{n \rightarrow \infty} f_n(x)$. Then $h(x) \leq f(x)$. Why? Now for each $\varepsilon > 0$, and fixed x , there exists z_n such that $f_n(x) + \varepsilon > f(z_n) + nd(z_n, x)$. Why? Therefore, $z_n \rightarrow x$. Why?

- (c) Then

$$\begin{aligned} h(x) + \varepsilon &= \lim_{n \rightarrow \infty} f_n(x) + \varepsilon \geq \lim_{n \rightarrow \infty} \inf (f(z_n) + nd(z_n, x)) \\ &\geq \lim_{n \rightarrow \infty} \inf f(z_n) \geq f(x) \end{aligned}$$

Why? Therefore, $h(x) \geq f(x)$ and so they are equal. Why?

- (d) Now consider $f : X \rightarrow (-\infty, \infty)$ and is lower semicontinuous as just explained. Consider $\frac{\pi}{2} + \arctan f(x) \equiv g(x)$. Then $\arctan f(x) \in (-\frac{\pi}{2}, \frac{\pi}{2})$ because f has real values. Then $g(x)$ is also lower semicontinuous having values in $(0, \pi)$. Why? By what was just shown, there exists $g_n(x) \uparrow g(x)$ where each g_n is continuous. Consider $f_n(x) \equiv \tan(g_n(x) - \frac{\pi}{2})$. Then f_n is continuous and increases to $f(x)$.

31. Generalize the above problem to the case where f is an upper semicontinuous real valued function. That is, $f(x) \geq \limsup_{n \rightarrow \infty} f(x_n)$ whenever $x_n \rightarrow x$. Show there are continuous functions $\{f_n(x)\}$ such that $f_n(x) \downarrow f(x)$. **Hint** To save trouble, maybe show that f is upper semicontinuous if and only if $-f$ is lower semicontinuous. Then maybe you could just use the above problem.
32. What if f is lower (upper) semicontinuous with values in $[-\infty, \infty]$? In this case, you consider $[-\infty, \infty]$ as a metric space as follows: $d(x, y) \equiv |\arctan(x) - \arctan(y)|$. Then you can generalize the above problems to show that if f is lower semicontinuous with values into $[-\infty, \infty]$ then it is the increasing limit of continuous functions with values in $[-\infty, \infty]$. Note that in this case a function identically equal to ∞ would be continuous so this is a rather odd sort of thing, a little different from what we normally like to consider. Check the details and explain why in this setting, the lower semicontinuous functions are exactly pointwise limits of increasing sequences of continuous functions and the upper semicontinuous functions are exactly pointwise limits of decreasing sequences of continuous functions.
33. This is a nice result in Taylor [57]. For a nonempty set T , ∂T is the set of points p such that $B(p, r)$ contains points of T and points of T^C for each $r > 0$. Suppose you have T a proper subset of a metric space and S is a connected, nonempty set such that $S \cap T \neq \emptyset, S \cap T^C \neq \emptyset$. Show that S must contain a point of ∂T .
34. Zorn's lemma is as follows: You have a nonempty partially ordered set \mathcal{F} with the partial order denoted by \prec and suppose you have the property that every totally ordered subset of \mathcal{F} has an upper bound. Show that it follows that there exists a maximal element $f \in \mathcal{F}$ such that if $f \prec g$ then $f = g$. **Hint:** Use the Hausdorff maximal theorem to show this. In fact, this is equivalent to the Hausdorff maximal theorem.

Chapter 4

Linear Spaces

The thing which is missing in the above material about metric spaces is any kind of algebra. In most applications, we are interested in adding things and multiplying things by scalars and so forth. This requires the notion of a vector space, also called a linear space. The simplest example is \mathbb{R}^n which is described next.

In this chapter, \mathbb{F} will refer to either \mathbb{R} or \mathbb{C} . It doesn't make any difference to the arguments which it is and so \mathbb{F} is written to symbolize whichever you wish to think about. When it is desired to emphasize that certain quantities are vectors, bold face will often be used. This is not necessarily done consistently. Sometimes context is considered sufficient.

4.1 Algebra in \mathbb{F}^n , Vector Spaces

There are exactly two algebraic operations done with elements of \mathbb{F}^n . One is addition and the other is multiplication by numbers, called scalars. In the case of \mathbb{C}^n the scalars are complex numbers while in the case of \mathbb{R}^n the only allowed scalars are real numbers. Thus, the scalars always come from \mathbb{F} in either case.

Definition 4.1.1 *If $\mathbf{x} \in \mathbb{F}^n$ and $a \in \mathbb{F}$, also called a scalar, then $a\mathbf{x} \in \mathbb{F}^n$ is defined by*

$$a\mathbf{x} = a(x_1, \dots, x_n) \equiv (ax_1, \dots, ax_n). \quad (4.1)$$

This is known as scalar multiplication. If $\mathbf{x}, \mathbf{y} \in \mathbb{F}^n$ then $\mathbf{x} + \mathbf{y} \in \mathbb{F}^n$ and is defined by

$$\begin{aligned} \mathbf{x} + \mathbf{y} &= (x_1, \dots, x_n) + (y_1, \dots, y_n) \\ &\equiv (x_1 + y_1, \dots, x_n + y_n) \end{aligned} \quad (4.2)$$

the points in \mathbb{F}^n are also referred to as vectors.

Actually, in dealing with vectors in \mathbb{F}^n , it is more customary in linear algebra to write them as column vectors. To save space, I will sometimes write $(x_1, \dots, x_n)^T$ to indicate the column vector having x_1 on the top and x_n on the bottom. With this definition, the algebraic properties satisfy the conclusions of the following theorem. These conclusions are called the vector space axioms. Any time you have a set and a field of scalars satisfying the axioms of the following theorem, it is called a vector space or linear space.

Theorem 4.1.2 *For $\mathbf{v}, \mathbf{w} \in \mathbb{F}^n$ and α, β scalars, (real numbers), the following hold.*

$$\mathbf{v} + \mathbf{w} = \mathbf{w} + \mathbf{v}, \quad (4.3)$$

the commutative law of addition,

$$(\mathbf{v} + \mathbf{w}) + \mathbf{z} = \mathbf{v} + (\mathbf{w} + \mathbf{z}), \quad (4.4)$$

the associative law for addition,

$$\mathbf{v} + \mathbf{0} = \mathbf{v}, \quad (4.5)$$

the existence of an additive identity,

$$\mathbf{v} + (-\mathbf{v}) = \mathbf{0}, \quad (4.6)$$

the existence of an additive inverse, Also

$$\alpha(v + w) = \alpha v + \alpha w, \quad (4.7)$$

$$(\alpha + \beta)v = \alpha v + \beta v, \quad (4.8)$$

$$\alpha(\beta v) = \alpha\beta(v), \quad (4.9)$$

$$1v = v. \quad (4.10)$$

In the above $\mathbf{0} = (0, \dots, 0)$.

You should verify these properties all hold. For example, consider 4.7

$$\begin{aligned} \alpha(v + w) &= \alpha(v_1 + w_1, \dots, v_n + w_n) \\ &= (\alpha(v_1 + w_1), \dots, \alpha(v_n + w_n)) \\ &= (\alpha v_1 + \alpha w_1, \dots, \alpha v_n + \alpha w_n) \\ &= (\alpha v_1, \dots, \alpha v_n) + (\alpha w_1, \dots, \alpha w_n) \\ &= \alpha v + \alpha w. \end{aligned}$$

As usual subtraction is defined as $x - y \equiv x + (-y)$.

4.2 Subspaces Spans and Bases

As mentioned above, \mathbb{F}^n is an example of a vector space. In dealing with vector spaces, the concept of linear combination is fundamental. When one considers only algebraic considerations, it makes no difference what field of scalars you are using. It could be \mathbb{R} , \mathbb{C} , \mathbb{Q} or even a field of residue classes. However, go ahead and think \mathbb{R} or \mathbb{C} since the subject of interest here is analysis.

Definition 4.2.1 Let $\{x_1, \dots, x_p\}$ be vectors in a vector space Y having the field of scalars \mathbb{F} . A linear combination is any expression of the form $\sum_{i=1}^p c_i x_i$ where the c_i are scalars. The set of all linear combinations of these vectors is called $\text{span}(x_1, \dots, x_p)$. A vector v is said to be in the span of some set S of vectors if v is a linear combination of vectors of S . **This means: finite linear combination.** If $V \subseteq Y$, then V is called a subspace if it contains $\mathbf{0}$ and whenever α, β are scalars and u and v are vectors of V , it follows $\alpha u + \beta v \in V$. That is, it is “closed under the algebraic operations of vector addition and scalar multiplication” and is therefore, a vector space. A linear combination of vectors is said to be trivial if all the scalars in the linear combination equal zero. A set of vectors is said to be linearly independent if the only linear combination of these vectors which equals the zero vector is the trivial linear combination. Thus $\{x_1, \dots, x_n\}$ is called linearly independent if whenever $\sum_{k=1}^n c_k x_k = \mathbf{0}$, it follows that all the scalars, c_k equal zero. A set of vectors, $\{x_1, \dots, x_n\}$, is called linearly dependent if it is not linearly independent. Thus the set of vectors is linearly dependent if there exist scalars, $c_i, i = 1, \dots, n$, not all zero such that $\sum_{k=1}^n c_k x_k = \mathbf{0}$.

Lemma 4.2.2 A set of vectors $\{x_1, \dots, x_n\}$ is linearly independent if and only if none of the vectors can be obtained as a linear combination of the others.

Proof: Suppose first that $\{x_1, \dots, x_n\}$ is linearly independent. If

$$x_k = \sum_{j \neq k} c_j x_j,$$

then $0 = 1x_k + \sum_{j \neq k} (-c_j)x_j$, a nontrivial linear combination, contrary to assumption. This shows that if the set is linearly independent, then none of the vectors is a linear combination of the others.

Now suppose no vector is a linear combination of the others. Is $\{x_1, \dots, x_n\}$ linearly independent? If it is not, there exist scalars, c_i , not all zero such that $\sum_{i=1}^n c_i x_i = 0$. Say $c_k \neq 0$. Then you can solve for x_k as $x_k = \sum_{j \neq k} (-c_j/c_k)x_j$ contrary to assumption. This proves the lemma. ■

The following is called the exchange theorem.

Theorem 4.2.3

$$\text{span}(u_1, \dots, u_r) \subseteq \text{span}(v_1, \dots, v_s) \equiv V$$

and $\{u_1, \dots, u_r\}$ are linearly independent, then $r \leq s$.

Proof: Suppose $r > s$. Let F_p denote the first p vectors in $\{u_1, \dots, u_r\}$. Let F_0 denote the empty set. Let E_p denote a finite list of vectors of $\{v_1, \dots, v_s\}$ and let $|E_p|$ denote the number of vectors in the list. Note that, by assumption, $\text{span}(F_0, E_s) = V$. For $0 \leq p \leq s$, let E_p have the property $\text{span}(F_p, E_p) = V$ and $|E_p|$ is as small as possible for this to happen. If $|E_p| = 0$, then $\text{span}(F_p) = V$ which would imply that, since $r > s \geq p$, $u_r \in \text{span}(F_s)$ contradicting the linear independence of $\{u_1, \dots, u_r\}$. Assume then that $|E_p| > 0$. Then $u_{p+1} \in \text{span}(F_p, E_p)$ and so there are constants, c_1, \dots, c_p and d_1, \dots, d_m such that $u_{p+1} = \sum_{i=1}^p c_i u_i + \sum_{j=1}^m d_j v_j$ for $\{z_1, \dots, z_m\} \subseteq \{v_1, \dots, v_s\}$. Then not all the d_i can equal zero because this would violate the linear independence of the $\{u_1, \dots, u_r\}$. Therefore, you can solve for one of the z_k as a linear combination of $\{u_1, \dots, u_{p+1}\}$ and the other z_j . Thus you can change F_p to F_{p+1} and include one fewer vector in E_{p+1} with $\text{span}(F_{p+1}, E_{p+1}) = V$ and so $|E_{p+1}| < |E_p|$ contrary to the claim that $|E_p|$ was as small as possible. Thus $|E_p| = 0$ after all and so a contradiction results.

Alternate proof: Recall from linear algebra that if you have A an $m \times n$ matrix where $m < n$ so there are more columns than rows, then there exists a nonzero solution x to the equation $Ax = 0$. Recall why this was. You must have free variables. Then by assumption, you have $u_j = \sum_{i=1}^s a_{ij} v_i$. If $s < r$, then the matrix (a_{ij}) has more columns than rows and so there exists a nonzero vector $x \in \mathbb{R}^r$ such that $\sum_{j=1}^r a_{ij} x_j = 0$. Then consider the following.

$$\sum_{j=1}^r x_j u_j = \sum_{j=1}^r x_j \sum_{i=1}^s a_{ij} v_i = \sum_i \sum_j a_{ij} x_j v_i = \sum_i 0 v_i = 0$$

and since not all $x_j = 0$, this contradicts the independence of $\{u_1, \dots, u_r\}$. ■

Definition 4.2.4 A finite set of vectors, $\{x_1, \dots, x_r\}$ is a basis for a vector space V if

$$\text{span}(x_1, \dots, x_r) = V$$

and $\{x_1, \dots, x_r\}$ is linearly independent. Thus if $v \in V$ there exist unique scalars, v_1, \dots, v_r such that $v = \sum_{i=1}^r v_i x_i$. These scalars are called the components of v with respect to the basis $\{x_1, \dots, x_r\}$ and $\{x_1, \dots, x_r\}$ are said to “span” V .

Corollary 4.2.5 Let $\{x_1, \dots, x_r\}$ and $\{y_1, \dots, y_s\}$ be two bases¹ of \mathbb{F}^n . Then $r = s = n$. More generally, if you have two bases for a vector space V then they have the same number of vectors.

Proof: From the exchange theorem, Theorem 4.2.3, if

$$\{x_1, \dots, x_r\}, \{y_1, \dots, y_s\}$$

are two bases for V , then $r \leq s$ and $s \leq r$. Now note the vectors,

$$e_i = \overbrace{(0, \dots, 0, 1, 0, \dots, 0)}^{1 \text{ is in the } i^{\text{th}} \text{ slot}}^T$$

for $i = 1, 2, \dots, n$ are a basis for \mathbb{F}^n . ■

Lemma 4.2.6 Let $\{v_1, \dots, v_r\}$ be a set of vectors. Then $V \equiv \text{span}(v_1, \dots, v_r)$ is a subspace.

Proof: Suppose α, β are two scalars and let $\sum_{k=1}^r c_k v_k$ and $\sum_{k=1}^r d_k v_k$ are two elements of V . What about $\alpha \sum_{k=1}^r c_k v_k + \beta \sum_{k=1}^r d_k v_k$? Is it also in V ?

$$\alpha \sum_{k=1}^r c_k v_k + \beta \sum_{k=1}^r d_k v_k = \sum_{k=1}^r (\alpha c_k + \beta d_k) v_k \in V$$

so the answer is yes. It is clear that 0 is in $\text{span}(v_1, \dots, v_r)$. This proves the lemma. ■

Definition 4.2.7 Let V be a vector space. It is finite dimensional when it has a basis of finitely many vectors. Otherwise, it is infinite dimensional. Then $\dim(V)$ read as the dimension of V is the number of vectors in a basis.

Of course you should wonder right now whether an arbitrary subspace of a finite dimensional vector space even has a basis. In fact it does and this is in the next theorem. First, here is an interesting lemma.

Lemma 4.2.8 Suppose $v \notin \text{span}(u_1, \dots, u_k)$ and $\{u_1, \dots, u_k\}$ is linearly independent. Then $\{u_1, \dots, u_k, v\}$ is also linearly independent.

Proof: Suppose $\sum_{i=1}^k c_i u_i + d v = 0$. It is required to verify that each $c_i = 0$ and that $d = 0$. But if $d \neq 0$, then you can solve for v as a linear combination of the vectors, $\{u_1, \dots, u_k\}$, $v = -\sum_{i=1}^k \left(\frac{c_i}{d}\right) u_i$ contrary to assumption. Therefore, $d = 0$. But then $\sum_{i=1}^k c_i u_i = 0$ and the linear independence of $\{u_1, \dots, u_k\}$ implies each $c_i = 0$ also. ■

Theorem 4.2.9 Let V be a nonzero subspace of Y a finite dimensional vector space having dimension n . Then V has a basis.

¹This is the plural form of basis. We could say basiss but it would involve an inordinate amount of hissing as in "The sixth shiek's sixth sheep is sick". This is the reason that bases is used instead of basiss.

Proof: Let $v_1 \in V$ where $v_1 \neq 0$. If $\text{span}\{v_1\} = V$, stop. $\{v_1\}$ is a basis for V . Otherwise, there exists $v_2 \in V$ which is not in $\text{span}\{v_1\}$. By Lemma 4.2.8 $\{v_1, v_2\}$ is a linearly independent set of vectors. If $\text{span}\{v_1, v_2\} = V$ stop, $\{v_1, v_2\}$ is a basis for V . If $\text{span}\{v_1, v_2\} \neq V$, then there exists $v_3 \notin \text{span}\{v_1, v_2\}$ and $\{v_1, v_2, v_3\}$ is a larger linearly independent set of vectors. Continuing this way, the process must stop before $n + 1$ steps because if not, it would be possible to obtain $n + 1$ linearly independent vectors contrary to the exchange theorem, Theorem 4.2.3, and the assumed dimension of V . ■

In words the following corollary states that any linearly independent set of vectors can be enlarged to form a basis.

Corollary 4.2.10 *Let V be a subspace of Y , a finite dimensional vector space of dimension n and let $\{v_1, \dots, v_r\}$ be a linearly independent set of vectors in V . Then either it is a basis for V or there exist vectors, v_{r+1}, \dots, v_s such that*

$$\{v_1, \dots, v_r, v_{r+1}, \dots, v_s\}$$

is a basis for V .

Proof: This follows immediately from the proof of Theorem 4.2.9. You do exactly the same argument except you start with $\{v_1, \dots, v_r\}$ rather than $\{v_1\}$. ■

It is also true that any spanning set of vectors can be restricted to obtain a basis.

Theorem 4.2.11 *Let V be a subspace of Y , a finite dimensional vector space of dimension n and suppose $\text{span}(u_1, \dots, u_p) = V$ where the u_i are nonzero vectors. Then there exist vectors, $\{v_1, \dots, v_r\}$ such that $\{v_1, \dots, v_r\} \subseteq \{u_1, \dots, u_p\}$ and $\{v_1, \dots, v_r\}$ is a basis for V .*

Proof: Let r be the smallest positive integer with the property that for some set,

$$\{v_1, \dots, v_r\} \subseteq \{u_1, \dots, u_p\}, \text{span}(v_1, \dots, v_r) = V.$$

Then $r \leq p$ and it must be the case that $\{v_1, \dots, v_r\}$ is linearly independent because if it were not so, one of the vectors, say v_k would be a linear combination of the others. But then you could delete this vector from $\{v_1, \dots, v_r\}$ and the resulting list of $r - 1$ vectors would still span V contrary to the definition of r . ■

4.3 Inner Product and Normed Linear Spaces

4.3.1 The Inner Product in \mathbb{F}^n

To do calculus, you must understand what you mean by distance. For functions of one variable, the distance was provided by the absolute value of the difference of two numbers. This must be generalized to \mathbb{F}^n and to more general situations.

Definition 4.3.1 *Let $x, y \in \mathbb{F}^n$. Thus $x = (x_1, \dots, x_n)$ where each $x_k \in \mathbb{F}$ and a similar formula holding for y . Then the inner product of these two vectors is defined to be*

$$(x, y) \equiv \sum_j x_j \bar{y}_j \equiv x_1 \bar{y}_1 + \dots + x_n \bar{y}_n.$$

Sometimes it is denoted as $x \cdot y$.

Notice how you put the conjugate on the entries of the vector \mathbf{y} . It makes no difference if the vectors happen to be real vectors but with complex vectors you must involve a conjugate. The reason for this is that when you take the inner product of a vector with itself, you want to get the square of the length of the vector, a positive number. Placing the conjugate on the components of \mathbf{y} in the above definition assures this will take place. Thus $(\mathbf{x}, \mathbf{x}) = \sum_j x_j \bar{x}_j = \sum_j |x_j|^2 \geq 0$. If you didn't place a conjugate as in the above definition, things wouldn't work out correctly. For example, $(1+i)^2 + 2^2 = 4 + 2i$ and this is not a positive number.

The following properties of the inner product follow immediately from the definition and you should verify each of them.

Properties of the inner product:

1. $(\mathbf{u}, \mathbf{v}) = \overline{(\mathbf{v}, \mathbf{u})}$
2. If a, b are numbers and $\mathbf{u}, \mathbf{v}, \mathbf{z}$ are vectors then $((a\mathbf{u} + b\mathbf{v}), \mathbf{z}) = a(\mathbf{u}, \mathbf{z}) + b(\mathbf{v}, \mathbf{z})$.
3. $(\mathbf{u}, \mathbf{u}) \geq 0$ and it equals 0 if and only if $\mathbf{u} = \mathbf{0}$.

Note this implies $(\mathbf{x}, \alpha\mathbf{y}) = \bar{\alpha}(\mathbf{x}, \mathbf{y})$ because

$$(\mathbf{x}, \alpha\mathbf{y}) = \overline{(\alpha\mathbf{y}, \mathbf{x})} = \overline{\alpha(\mathbf{y}, \mathbf{x})} = \bar{\alpha}(\mathbf{x}, \mathbf{y})$$

The norm is defined as follows.

Definition 4.3.2 For $\mathbf{x} \in \mathbb{F}^n$, $|\mathbf{x}| \equiv \left(\sum_{k=1}^n |x_k|^2 \right)^{1/2} = (\mathbf{x}, \mathbf{x})^{1/2}$.

4.3.2 General Inner Product Spaces

Any time you have a vector space which possesses an inner product, something satisfying the properties 1 - 3 above, it is called an inner product space.

Here is a fundamental inequality called the **Cauchy Schwarz inequality** which holds in any inner product space. First here is a simple lemma.

Lemma 4.3.3 If $z \in \mathbb{F}$ there exists $\theta \in \mathbb{F}$ such that $\theta z = |z|$ and $|\theta| = 1$.

Proof: Let $\theta = 1$ if $z = 0$ and otherwise, let $\theta = \frac{\bar{z}}{|z|}$. Recall that for $z = x + iy$, $\bar{z} = x - iy$

and $\bar{z}z = |z|^2$. In case z is real, there is no change in the above. ■

Theorem 4.3.4 (Cauchy Schwarz) Let H be an inner product space. The following inequality holds for \mathbf{x} and $\mathbf{y} \in H$.

$$|(\mathbf{x}, \mathbf{y})| \leq (\mathbf{x}, \mathbf{x})^{1/2} (\mathbf{y}, \mathbf{y})^{1/2} \quad (4.11)$$

Equality holds in this inequality if and only if one vector is a multiple of the other.

Proof: Let $\theta \in \mathbb{F}$ such that $|\theta| = 1$ and $\theta(\mathbf{x}, \mathbf{y}) = |(\mathbf{x}, \mathbf{y})|$. Consider

$$p(t) \equiv (\mathbf{x} + \bar{\theta}t\mathbf{y}, \mathbf{x} + t\bar{\theta}\mathbf{y})$$

where $t \in \mathbb{R}$. Then from the above list of properties of the inner product,

$$\begin{aligned}
 0 &\leq p(t) = (x, x) + t\theta(x, y) + t\bar{\theta}(y, x) + t^2(y, y) \\
 &= (x, x) + t\theta(x, y) + t\overline{\theta(x, y)} + t^2(y, y) \\
 &= (x, x) + 2t\operatorname{Re}(\theta(x, y)) + t^2(y, y) \\
 &= (x, x) + 2t|(x, y)| + t^2(y, y)
 \end{aligned} \tag{4.12}$$

and this must hold for all $t \in \mathbb{R}$. Therefore, if $(y, y) = 0$ it must be the case that $|(x, y)| = 0$ also since otherwise the above inequality would be violated. Therefore, in this case, $|(x, y)| \leq (x, x)^{1/2}(y, y)^{1/2}$. On the other hand, if $(y, y) \neq 0$, then $p(t) \geq 0$ for all t means the graph of $y = p(t)$ is a parabola which opens up and it either has exactly one real zero in the case its vertex touches the t axis or it has no real zeros. From the quadratic formula this happens exactly when $4|(x, y)|^2 - 4(x, x)(y, y) \leq 0$ which is equivalent to 4.11.

It is clear from a computation that if one vector is a scalar multiple of the other that equality holds in 4.11. Conversely, suppose equality does hold. Then this is equivalent to saying $4|(x, y)|^2 - 4(x, x)(y, y) = 0$ and so from the quadratic formula, there exists one real zero to $p(t) = 0$. Call it t_0 . Then

$$p(t_0) = (x + \bar{\theta}t_0y, x + t_0\bar{\theta}y) = |x + \bar{\theta}t_0y|^2 = 0$$

and so $x = -\bar{\theta}t_0y$. ■

Note that in establishing the inequality, I only used part of the above properties of the inner product. It was not necessary to use the one which says that if $(x, x) = 0$ then $x = 0$. That was only used to consider the case of equality.

Now the length of a vector can be defined.

Definition 4.3.5 Let $z \in H$. Then $|z| \equiv (z, z)^{1/2}$.

Theorem 4.3.6 For length defined in Definition 4.3.5, the following hold.

$$|z| \geq 0 \text{ and } |z| = 0 \text{ if and only if } z = 0 \tag{4.13}$$

$$\text{If } \alpha \text{ is a scalar, } |\alpha z| = |\alpha||z| \tag{4.14}$$

$$|z + w| \leq |z| + |w|. \tag{4.15}$$

Proof: The first two claims are left as exercises. To establish the third,

$$\begin{aligned}
 |z + w|^2 &\equiv (z + w, z + w) \\
 &= (z, z) + (w, w) + (w, z) + (z, w) \\
 &= |z|^2 + |w|^2 + 2\operatorname{Re}(w, z) \\
 &\leq |z|^2 + |w|^2 + 2|(w, z)| \\
 &\leq |z|^2 + |w|^2 + 2|w||z| = (|z| + |w|)^2.
 \end{aligned}$$

Note that in an inner product space, you can define $d(x, y) \equiv |x - y|$ and this is a metric for this inner product space. This follows from the above since d satisfies the conditions for a metric,

$$d(x, y) = d(y, x), \quad d(x, y) \geq 0 \text{ and equals } 0 \text{ if and only if } x = y$$

$$d(x, y) + d(y, z) = |x - y| + |y - z| \geq |x - y + y - z| = |x - z| = d(x, z).$$

It follows that all the theory of metric spaces developed earlier applies to this situation.

4.3.3 Normed Vector Spaces

The best sort of a norm is one which comes from an inner product. However, any vector space V which has a function $\|\cdot\|$ which maps V to $[0, \infty)$ is called a normed vector space if $\|\cdot\|$ satisfies 4.13 - 4.15. That is

$$\|z\| \geq 0 \text{ and } \|z\| = 0 \text{ if and only if } z = 0 \quad (4.16)$$

$$\text{If } \alpha \text{ is a scalar, } \|\alpha z\| = |\alpha| \|z\| \quad (4.17)$$

$$\|z + w\| \leq \|z\| + \|w\|. \quad (4.18)$$

The last inequality above is called the triangle inequality. Another version of this is

$$|\|z\| - \|w\|| \leq \|z - w\| \quad (4.19)$$

To see that 4.19 holds, note $\|z\| = \|z - w + w\| \leq \|z - w\| + \|w\|$ which implies $\|z\| - \|w\| \leq \|z - w\|$ and now switching z and w , yields $\|w\| - \|z\| \leq \|z - w\|$ which implies 4.19.

Any normed vector space is a metric space, the distance given by $d(x, y) \equiv \|x - y\|$. This satisfies all the axioms of a distance. Therefore, any normed linear space is a metric space with this metric and all the theory of metric spaces applies.

Definition 4.3.7 When X is a normed linear space which is also complete, it is called a Banach space.

A Banach space may or may not be finite dimensional but it is always a linear space or vector space. The field of scalars will always be \mathbb{R} or \mathbb{C} at least in this book. More is said about Banach spaces later.

4.3.4 The p Norms

Examples of norms are the p norms on \mathbb{C}^n for $p \neq 2$. These do not come from an inner product but they are norms just the same.

Definition 4.3.8 Let $x \in \mathbb{C}^n$. Then define for $p \geq 1$,

$$\|x\|_p \equiv \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}.$$

The following inequality is called Holder's inequality.

Proposition 4.3.9 For $x, y \in \mathbb{C}^n$,

$$\sum_{i=1}^n |x_i| |y_i| \leq \left(\sum_{i=1}^n |x_i|^p \right)^{1/p} \left(\sum_{i=1}^n |y_i|^{p'} \right)^{1/p'}$$

The proof will depend on the following lemma shown later.

Lemma 4.3.10 If $a, b \geq 0$ and p' is defined by $\frac{1}{p} + \frac{1}{p'} = 1$, then

$$ab \leq \frac{a^p}{p} + \frac{b^{p'}}{p'}.$$

Proof of the Proposition: If \mathbf{x} or \mathbf{y} equals the zero vector there is nothing to prove. Therefore, assume they are both nonzero. Let $A = (\sum_{i=1}^n |x_i|^p)^{1/p}$ and $B = (\sum_{i=1}^n |y_i|^{p'})^{1/p'}$. Then using Lemma 4.3.10,

$$\begin{aligned} \sum_{i=1}^n \frac{|x_i|}{A} \frac{|y_i|}{B} &\leq \sum_{i=1}^n \left[\frac{1}{p} \left(\frac{|x_i|}{A} \right)^p + \frac{1}{p'} \left(\frac{|y_i|}{B} \right)^{p'} \right] \\ &= \frac{1}{p} \frac{1}{A^p} \sum_{i=1}^n |x_i|^p + \frac{1}{p'} \frac{1}{B^{p'}} \sum_{i=1}^n |y_i|^{p'} \\ &= \frac{1}{p} + \frac{1}{p'} = 1 \end{aligned}$$

and so $\sum_{i=1}^n |x_i| |y_i| \leq AB = (\sum_{i=1}^n |x_i|^p)^{1/p} (\sum_{i=1}^n |y_i|^{p'})^{1/p'}$. ■

Theorem 4.3.11 *The p norms do indeed satisfy the axioms of a norm.*

Proof: It is obvious that $\|\cdot\|_p$ does indeed satisfy most of the norm axioms. The only one that is not clear is the triangle inequality. To save notation write $\|\cdot\|$ in place of $\|\cdot\|_p$ in what follows. Note also that $\frac{p}{p'} = p - 1$. Then using the Holder inequality,

$$\begin{aligned} \|\mathbf{x} + \mathbf{y}\|^p &= \sum_{i=1}^n |x_i + y_i|^p \leq \sum_{i=1}^n |x_i + y_i|^{p-1} |x_i| + \sum_{i=1}^n |x_i + y_i|^{p-1} |y_i| \\ &= \sum_{i=1}^n |x_i + y_i|^{\frac{p}{p'}} |x_i| + \sum_{i=1}^n |x_i + y_i|^{\frac{p}{p'}} |y_i| \\ &\leq \left(\sum_{i=1}^n |x_i + y_i|^p \right)^{1/p'} \left[\left(\sum_{i=1}^n |x_i|^p \right)^{1/p} + \left(\sum_{i=1}^n |y_i|^p \right)^{1/p} \right] \\ &= \|\mathbf{x} + \mathbf{y}\|^{p/p'} (\|\mathbf{x}\|_p + \|\mathbf{y}\|_p) \end{aligned}$$

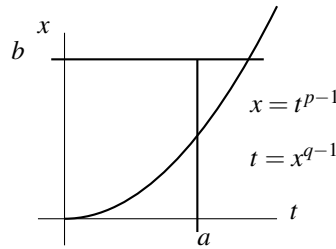
so dividing by $\|\mathbf{x} + \mathbf{y}\|^{p/p'}$, it follows

$$\|\mathbf{x} + \mathbf{y}\|^p \|\mathbf{x} + \mathbf{y}\|^{-p/p'} = \|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\|_p + \|\mathbf{y}\|_p$$

$\left(p - \frac{p}{p'} = p \left(1 - \frac{1}{p'} \right) = p \frac{1}{p} = 1 \right)$. ■

It only remains to prove Lemma 4.3.10.

Proof of the lemma: Let $p' = q$ to save on notation and consider the following picture:



$$ab \leq \int_0^a t^{p-1} dt + \int_0^b x^{q-1} dx = \frac{a^p}{p} + \frac{b^q}{q}.$$

Note equality occurs when $a^p = b^q$. ■

Alternate proof of the lemma: First note that if either a or b are zero, then there is nothing to show so we can assume $b, a > 0$. Let $b > 0$ and let

$$f(a) = \frac{a^p}{p} + \frac{b^q}{q} - ab$$

Then the second derivative of f is positive on $(0, \infty)$ so its graph is convex. Also $f(0) > 0$ and $\lim_{a \rightarrow \infty} f(a) = \infty$. Then a short computation shows that there is only one critical point, where f is minimized and this happens when a is such that $a^p = b^q$. At this point,

$$f(a) = b^q - b^{q/p}b = b^q - b^{q-1}b = 0$$

Therefore, $f(a) \geq 0$ for all a and this proves the lemma. ■

Another example of a very useful norm on \mathbb{F}^n is the norm $\|\cdot\|_\infty$ defined by

$$\|\mathbf{x}\|_\infty \equiv \max \{|x_k| : k = 1, 2, \dots, n\}$$

You should verify that this satisfies all the axioms of a norm. Here is the triangle inequality.

$$\begin{aligned} \|\mathbf{x} + \mathbf{y}\|_\infty &= \max_k \{|x_k + y_k|\} \leq \max_k \{|x_k| + |y_k|\} \\ &\leq \max_k \{|x_k|\} + \max_k \{|y_k|\} = \|\mathbf{x}\|_\infty + \|\mathbf{y}\|_\infty \end{aligned}$$

It turns out that in terms of analysis, it makes **absolutely no difference** which norm you use. This will be explained later. First is a short review of the notion of orthonormal bases which is not needed directly in what follows but is sufficiently important to include.

4.3.5 Orthonormal Bases

Not all bases for an inner product space H are created equal. The best bases are orthonormal.

Definition 4.3.12 Suppose $\{v_1, \dots, v_k\}$ is a set of vectors in an inner product space H . It is an orthonormal set if

$$(v_i, v_j) = \delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

Every orthonormal set of vectors is automatically linearly independent.

Proposition 4.3.13 Suppose $\{v_1, \dots, v_k\}$ is an orthonormal set of vectors. Then it is linearly independent.

Proof: Suppose $\sum_{i=1}^k c_i v_i = \mathbf{0}$. Then taking inner products with

$$v_j, 0 = (\mathbf{0}, v_j) = \sum_i c_i (v_i, v_j) = \sum_i c_i \delta_{ij} = c_j.$$

Since j is arbitrary, this shows the set is linearly independent as claimed. ■

It turns out that if X is any subspace of H , then there exists an orthonormal basis for X . The process by which this is done is called the Gram Schmidt process.

Lemma 4.3.14 *Let X be a subspace of dimension n which is contained in an inner product space H . Let a basis for X be $\{x_1, \dots, x_n\}$. Then there exists an orthonormal basis for X , $\{u_1, \dots, u_n\}$ which has the property that for each $k \leq n$, $\text{span}(x_1, \dots, x_k) = \text{span}(u_1, \dots, u_k)$.*

Proof: Let $\{x_1, \dots, x_n\}$ be a basis for X . Let $u_1 \equiv x_1/|x_1|$. Thus for $k = 1$,

$$\text{span}(u_1) = \text{span}(x_1)$$

and $\{u_1\}$ is an orthonormal set. Now suppose for some $k < n$, u_1, \dots, u_k have been chosen such that $(u_j, u_l) = \delta_{jl}$ and $\text{span}(x_1, \dots, x_k) = \text{span}(u_1, \dots, u_k)$. Then define

$$u_{k+1} \equiv \frac{x_{k+1} - \sum_{j=1}^k (x_{k+1}, u_j) u_j}{\left| x_{k+1} - \sum_{j=1}^k (x_{k+1}, u_j) u_j \right|}, \quad (4.20)$$

where the denominator is not equal to zero because the x_j form a basis and so

$$x_{k+1} \notin \text{span}(x_1, \dots, x_k) = \text{span}(u_1, \dots, u_k)$$

Thus by induction,

$$u_{k+1} \in \text{span}(u_1, \dots, u_k, x_{k+1}) = \text{span}(x_1, \dots, x_k, x_{k+1}).$$

Also, $x_{k+1} \in \text{span}(u_1, \dots, u_k, u_{k+1})$ which is seen easily by solving 4.20 for x_{k+1} and it follows

$$\text{span}(x_1, \dots, x_k, x_{k+1}) = \text{span}(u_1, \dots, u_k, u_{k+1}).$$

If $l \leq k$, then denoting by C the scalar $\left| x_{k+1} - \sum_{j=1}^k (x_{k+1}, u_j) u_j \right|^{-1}$,

$$\begin{aligned} (u_{k+1}, u_l) &= C \left((x_{k+1}, u_l) - \sum_{j=1}^k (x_{k+1}, u_j) (u_j, u_l) \right) \\ &= C \left((x_{k+1}, u_l) - \sum_{j=1}^k (x_{k+1}, u_j) \delta_{lj} \right) \\ &= C((x_{k+1}, u_l) - (x_{k+1}, u_l)) = 0. \end{aligned}$$

The vectors, $\{u_j\}_{j=1}^n$, generated in this way are therefore an orthonormal basis because each vector has unit length. ■

The process by which these vectors were generated is called the Gram Schmidt process.

4.4 Equivalence of Norms

As mentioned above, it makes absolutely no difference which norm you decide to use. This holds in general finite dimensional normed spaces. First are some simple lemmas featuring one dimensional considerations. In this case, the distance is given by $d(x, y) = |x - y|$ and so the open balls are sets of the form $(x - \delta, x + \delta)$.

Also recall the Lemma 3.5.9 which is stated next for convenience.

Lemma 4.4.1 *The closed interval $[a, b]$ is compact.*

Corollary 4.4.2 *The set $Q \equiv [a, b] + i[c, d] \subseteq \mathbb{C}$ is compact, meaning*

$$\{x + iy : x \in [a, b], y \in [c, d]\}$$

Proof: Let $\{x_n + iy_n\}$ be a sequence in Q . Then there is a subsequence such that $\lim_{k \rightarrow \infty} x_{n_k} = x \in [a, b]$. There is a further subsequence such that $\lim_{l \rightarrow \infty} y_{n_{k_l}} = y \in [c, d]$. Thus, also $\lim_{l \rightarrow \infty} x_{n_{k_l}} = x$ because subsequences of convergent sequences converge to the same point. Therefore, from the way we measure the distance in \mathbb{C} , it follows that $\lim_{l \rightarrow \infty} (x_{n_{k_l}} + iy_{n_{k_l}}) = x + iy \in Q$. ■

The next corollary gives the definition of a closed disk and shows that, like a closed interval, a closed disk is compact.

Corollary 4.4.3 *In \mathbb{C} , let $D(z, r) \equiv \{w \in \mathbb{C} : |z - w| \leq r\}$. Then $D(z, r)$ is compact.*

Proof: Note that

$$D(z, r) \subseteq [\operatorname{Re} z - r, \operatorname{Re} z + r] + i[\operatorname{Im} z - r, \operatorname{Im} z + r]$$

which was just shown to be compact. Also, if $w_k \rightarrow w$ where $w_k \in D(z, r)$, then by the triangle inequality,

$$|z - w| = \lim_{k \rightarrow \infty} |z - w_k| \leq r$$

and so $D(z, r)$ is a closed subset of a compact set. Hence it is compact by Proposition 3.5.2. ■

Recall that sequentially compact and compact are the same in any metric space which is the context of the assertions here.

Lemma 4.4.4 *Let K_i be a nonempty compact set in \mathbb{F} . Then $P \equiv \prod_{i=1}^n K_i$ is compact in \mathbb{F}^n .*

Proof: Let $\{\mathbf{x}_k\}$ be a sequence in P . Taking a succession of subsequences as in the proof of Corollary 4.4.2, there exists a subsequence, still denoted as $\{\mathbf{x}_k\}$ such that if x_k^i is the i^{th} component of \mathbf{x}_k , then $\lim_{k \rightarrow \infty} x_k^i = x^i \in K_i$. Thus if \mathbf{x} is the vector of P whose i^{th} component is x^i ,

$$\lim_{k \rightarrow \infty} |\mathbf{x}_k - \mathbf{x}| \equiv \lim_{k \rightarrow \infty} \left(\sum_{i=1}^n |x_k^i - x^i|^2 \right)^{1/2} = 0$$

It follows that P is sequentially compact, hence compact. ■

A set K in \mathbb{F}^n is said to be bounded if it is contained in some ball $B(\mathbf{0}, r)$.

Theorem 4.4.5 *A set $K \subseteq \mathbb{F}^n$ is compact if it is closed and bounded. If $f : K \rightarrow \mathbb{R}$, then f achieves its maximum and its minimum on K .*

Proof: Say K is closed and bounded, being contained in $B(\mathbf{0}, r)$. Then if $\mathbf{x} \in K$, $|x_i| < r$ where x_i is the i^{th} component. Hence $K \subseteq \prod_{i=1}^n D(0, r)$, a compact set by Lemma 4.4.4. By Proposition 3.5.2, since K is a closed subset of a compact set, it is compact. The last claim is just the extreme value theorem, Theorem 3.7.2. ■

Definition 4.4.6 Let $\{v_1, \dots, v_n\}$ be a basis for V where $(V, \|\cdot\|)$ is a finite dimensional normed vector space with field of scalars equal to either \mathbb{R} or \mathbb{C} . Define $\theta : V \rightarrow \mathbb{F}^n$ as follows.

$$\theta \left(\sum_{j=1}^n \alpha_j v_j \right) \equiv \alpha \equiv (\alpha_1, \dots, \alpha_n)^T$$

Thus θ maps a vector to its coordinates taken with respect to a given basis.

The following fundamental lemma comes from the extreme value theorem for continuous functions defined on a compact set. Let

$$f(\alpha) \equiv \left\| \sum_i \alpha_i v_i \right\| \equiv \|\theta^{-1} \alpha\|$$

Then it is clear that f is a continuous function defined on \mathbb{F}^n . This is because $\alpha \rightarrow \sum_i \alpha_i v_i$ is a continuous map into V and from the triangle inequality $x \rightarrow \|x\|$ is continuous as a map from V to \mathbb{R} .

Lemma 4.4.7 There exists $\delta > 0$ and $\Delta \geq \delta$ such that

$$\delta = \min \{f(\alpha) : |\alpha| = 1\}, \quad \Delta = \max \{f(\alpha) : |\alpha| = 1\}$$

Also,

$$\delta |\alpha| \leq \|\theta^{-1} \alpha\| \leq \Delta |\alpha| \quad (4.21)$$

$$\delta \|\theta v\| \leq \|v\| \leq \Delta \|\theta v\| \quad (4.22)$$

Proof: These numbers exist thanks to Theorem 4.4.5. It cannot be that $\delta = 0$ because if it were, you would have $|\alpha| = 1$ but $\sum_{j=1}^n \alpha_j v_j = \mathbf{0}$ which is impossible since $\{v_1, \dots, v_n\}$ is linearly independent. The first of the above inequalities follows from $\delta \leq \left\| \theta^{-1} \frac{\alpha}{|\alpha|} \right\| = f\left(\frac{\alpha}{|\alpha|}\right) \leq \Delta$. The second follows from observing that $\theta^{-1} \alpha$ is a generic vector v in V . ■

Note that these inequalities yield the fact that convergence of the coordinates with respect to a given basis is equivalent to convergence of the vectors. More precisely, to say that $\lim_{k \rightarrow \infty} v^k = v$ is the same as saying that $\lim_{k \rightarrow \infty} \theta v^k = \theta v$. Indeed,

$$\delta \|\theta v_n - \theta v\| \leq \|v_n - v\| \leq \Delta \|\theta v_n - \theta v\|$$

Now we can draw several conclusions about $(V, \|\cdot\|)$ for V finite dimensional.

Theorem 4.4.8 Let $(V, \|\cdot\|)$ be a finite dimensional normed linear space. Then the compact sets are exactly those which are closed and bounded. Also $(V, \|\cdot\|)$ is complete. If K is a closed and bounded set in $(V, \|\cdot\|)$ and $f : K \rightarrow \mathbb{R}$, then f achieves its maximum and minimum on K .

Proof: First note that the inequalities 4.21 and 4.22 show that both θ^{-1} and θ are continuous. Thus these take convergent sequences to convergent sequences.

Let $\{w_k\}_{k=1}^\infty$ be a Cauchy sequence. Then from 4.22, $\{\theta w_k\}_{k=1}^\infty$ is a Cauchy sequence. Thanks to Theorem 4.4.5, it converges to some $\beta \in \mathbb{F}^n$. It follows that $\lim_{k \rightarrow \infty} \theta^{-1} \theta w_k = \lim_{k \rightarrow \infty} w_k = \theta^{-1} \beta \in V$. This shows completeness.

Next let K be a closed and bounded set. Let $\{w_k\} \subseteq K$. Then $\{\theta w_k\} \subseteq \theta K$ which is also a closed and bounded set thanks to the inequalities 4.21 and 4.22. Thus there is a subsequence still denoted with k such that $\theta w_k \rightarrow \beta \in \mathbb{F}^n$. Then as just done, $w_k \rightarrow \theta^{-1}\beta$. Since K is closed, it follows that $\theta^{-1}\beta \in K$.

This has just shown that a closed and bounded set in V is sequentially compact hence compact.

Finally, why are the only compact sets those which are closed and bounded? Let K be compact. If it is not bounded, then there is a sequence of points of K , $\{k^m\}_{m=1}^\infty$ such that $\|k^m\| \geq \|k^{m-1}\| + 1$. It follows that it cannot have a convergent subsequence because the points are further apart from each other than $1/2$. Indeed,

$$\|k^m - k^{m+1}\| \geq \|k^{m+1}\| - \|k^m\| \geq 1 > 1/2$$

Hence K is not sequentially compact and consequently it is not compact. It follows that K is bounded. If K is not closed, then there exists a limit point k which is not in K . (Recall that closed means it has all its limit points.) By Theorem 3.1.8, there is a sequence of distinct points having no repeats and none equal to k denoted as $\{k^m\}_{m=1}^\infty$ such that $k^m \rightarrow k$. Then this sequence $\{k^m\}$ fails to have a subsequence which converges to a point of K . Hence K is not sequentially compact. Thus, if K is compact then it is closed and bounded.

The last part is the extreme value theorem, Theorem 3.7.2. ■

Next is the theorem which states that any two norms on a finite dimensional vector space are equivalent.

Theorem 4.4.9 *Let $\|\cdot\|, \|\cdot\|_1$ be two norms on V a finite dimensional vector space. Then they are equivalent, which means there are constants $0 < a < b$ such that for all v ,*

$$a\|v\| \leq \|v\|_1 \leq b\|v\|$$

Proof: In Lemma 4.4.7, let δ, Δ go with $\|\cdot\|$ and $\hat{\delta}, \hat{\Delta}$ go with $\|\cdot\|_1$. Then using the inequalities of this lemma,

$$\|v\| \leq \Delta|\theta v| \leq \frac{\Delta}{\hat{\delta}}\|v\|_1 \leq \frac{\Delta\hat{\Delta}}{\hat{\delta}}|\theta v| \leq \frac{\Delta\hat{\Delta}}{\hat{\delta}}\|v\|$$

and so $\frac{\hat{\delta}}{\Delta}\|v\| \leq \|v\|_1 \leq \frac{\hat{\Delta}}{\hat{\delta}}\|v\|$. Thus the norms are equivalent. ■

It follows right away that the closed and open sets are the same with two different norms. Also, all considerations involving limits are unchanged from one norm to another.

Corollary 4.4.10 *Consider the metric spaces $(V, \|\cdot\|_1), (V, \|\cdot\|_2)$ where V has dimension n . Then a set is closed or open in one of these if and only if it is respectively closed or open in the other. In other words, the two metric spaces have exactly the same open and closed sets. Also, a set is bounded in one metric space if and only if it is bounded in the other.*

Proof: This follows from Theorem 3.6.2, the theorem about the equivalent formulations of continuity. Using this theorem, it follows from Theorem 4.4.9 that the identity map $I(x) \equiv x$ is continuous. The reason for this is that the inequality of this theorem implies that if $\|v^m - v\|_1 \rightarrow 0$ then $\|Iv^m - Iv\|_2 = \|I(v^m - v)\|_2 \rightarrow 0$ and the same holds on switching 1 and 2 in what was just written.

Therefore, the identity map takes open sets to open sets and closed sets to closed sets. In other words, the two metric spaces have the same open sets and the same closed sets.

Suppose S is bounded in $(V, \|\cdot\|_1)$. This means it is contained in $B(\mathbf{0}, r)_1$ where the subscript of 1 indicates the norm is $\|\cdot\|_1$. Let $\delta \|\cdot\|_1 \leq \|\cdot\|_2 \leq \Delta \|\cdot\|_1$ as described above. Then $S \subseteq B(\mathbf{0}, r)_1 \subseteq B(\mathbf{0}, \Delta r)_2$ so S is also bounded in $(V, \|\cdot\|_2)$. Similarly, if S is bounded in $\|\cdot\|_2$ then it is bounded in $\|\cdot\|_1$. ■

One can show that in the case of \mathbb{R} where it makes sense to consider sup and inf, convergence of Cauchy sequences can be shown to imply the other definition of completeness involving sup, and inf.

4.5 Covering Theorems

These covering theorems make sense on any finite dimensional normed linear space. There are two which are commonly used, the Vitali theorem and the Besicovitch theorem. The first adjusts the size of balls and the second does not. Of the two, it is the Besicovitch theorem which I will emphasize. However, the Vitali theorem is used more often and may be a little easier. I decided to place these theorems early in the book to emphasize that they only require a finite dimensional normed linear space.

4.5.1 Vitali Covering Theorem

The Vitali covering theorem is a profound result about coverings of a set in $(X, \|\cdot\|)$ with balls. Usually we are interested in \mathbb{R}^p with some norm. We will tacitly assume all balls have positive radius. They will not be single points. Before beginning the proof, here is a useful lemma.

Lemma 4.5.1 *In a normed linear space, $\overline{B(\mathbf{x}, r)} = \{\mathbf{y} : \|\mathbf{y} - \mathbf{x}\| \leq r\}$.*

Proof: It is clear that $\overline{B(\mathbf{x}, r)} \subseteq \{\mathbf{y} : \|\mathbf{y} - \mathbf{x}\| \leq r\}$ because if $\mathbf{y} \in \overline{B(\mathbf{x}, r)}$, then there exists a sequence of points of $B(\mathbf{x}, r)$, $\{\mathbf{x}_n\}$ such that $\|\mathbf{x}_n - \mathbf{y}\| \rightarrow 0$, $\|\mathbf{x}_n\| < r$. However, this requires that $\|\mathbf{x}_n\| \rightarrow \|\mathbf{y}\|$ and so $\|\mathbf{y}\| \leq r$. Now let \mathbf{y} be in the right side. It suffices to consider $\|\mathbf{y} - \mathbf{x}\| = 1$. Then you could consider for $t \in (0, 1)$, $\mathbf{x} + t(\mathbf{y} - \mathbf{x}) = \mathbf{z}(t)$. Then $\|\mathbf{z}(t) - \mathbf{x}\| = t\|\mathbf{y} - \mathbf{x}\| = tr < r$ and so $\mathbf{z}(t) \in B(\mathbf{x}, r)$. But also, $\|\mathbf{z}(t) - \mathbf{y}\| = (1-t)\|\mathbf{y} - \mathbf{x}\| = (1-t)r$ so $\lim_{t \rightarrow 0} \|\mathbf{z}(t) - \mathbf{y}\| = 0$ showing that $\mathbf{y} \in \overline{B(\mathbf{x}, r)}$. ■

Thus the usual way we think about the closure of a ball is completely correct in a normed linear space. Its limit points not in the ball are exactly \mathbf{y} such that $\|\mathbf{y} - \mathbf{x}\| = r$. Recall that this lemma is not always true in the context of a metric space. Recall the discrete metric for example, in which the distance between different points is 1 and distance between a point and itself is 0. In what follows I will use the result of this lemma without comment. Balls will be either open, closed or neither. I am going to use the Hausdorff maximal theorem, Theorem 2.8.2 because it yields a very simple argument. It can be done other ways however. In the argument, the balls are not necessarily open nor closed. \mathbf{y} is in $B(\mathbf{x}, r)$ will mean that $\|\mathbf{y} - \mathbf{x}\| < r$ or $\|\mathbf{y} - \mathbf{x}\| = r$.

Lemma 4.5.2 *Let \mathcal{F} be a nonempty collection of balls satisfying*

$$\infty > M \equiv \sup\{r : B(\mathbf{p}, r) \in \mathcal{F}\} > 0$$

and let $k \in (0, M)$. Then there exists $\mathcal{G} \subseteq \mathcal{F}$ such that

$$\text{If } B(\mathbf{p}, r) \in \mathcal{G}, \text{ then } r > k, \quad (4.23)$$

$$\text{If } B_1, B_2 \in \mathcal{G} \text{ then } \overline{B_1} \cap \overline{B_2} = \emptyset, \quad (4.24)$$

$$\mathcal{G} \text{ is maximal with respect to 4.23 and 4.24.} \quad (4.25)$$

By this is meant that if \mathcal{H} is a collection of balls satisfying 4.23 and 4.24, then \mathcal{H} cannot properly contain \mathcal{G} .

Proof: Let \mathfrak{S} denote a subset of \mathcal{F} such that 4.23 and 4.24 are satisfied. Since $k < M$, 4.23 is satisfied for some ball of \mathfrak{S} . Thus $\mathfrak{S} \neq \emptyset$. Partially order \mathfrak{S} with respect to set inclusion. Thus $\mathcal{A} \prec \mathcal{B}$ for \mathcal{A}, \mathcal{B} in \mathfrak{S} means that $\mathcal{A} \subseteq \mathcal{B}$. By the Hausdorff maximal theorem, there is a maximal chain in \mathfrak{S} denoted by \mathcal{C} . Then let \mathcal{G} be $\cup \mathcal{C}$. If B_1, B_2 are in \mathcal{C} , then since \mathcal{C} is a chain, both B_1, B_2 are in some element of \mathcal{C} and so $\overline{B_1} \cap \overline{B_2} = \emptyset$. The maximality of \mathcal{C} is violated if there is any other element of \mathfrak{S} which properly contains \mathcal{G} . ■

Proposition 4.5.3 Let \mathcal{F} be a collection of balls, and let

$$A \equiv \cup \{B : B \in \mathcal{F}\}.$$

Suppose

$$\infty > M \equiv \sup \{r : B(\mathbf{p}, r) \in \mathcal{F}\} > 0.$$

Then there exists $\mathcal{G} \subseteq \mathcal{F}$ such that \mathcal{G} consists of balls whose closures are disjoint and

$$A \subseteq \cup \{\widehat{B} : B \in \mathcal{G}\}$$

where for $B = B(\mathbf{x}, r)$ a ball, \widehat{B} denotes the open ball $B(\mathbf{x}, 5r)$.

Proof: Let \mathcal{G}_1 satisfy 4.23 - 4.25 for $k = \frac{2M}{3}$.

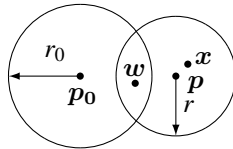
Suppose $\mathcal{G}_1, \dots, \mathcal{G}_{m-1}$ have been chosen for $m \geq 2$. Let $\overline{\mathcal{G}_i}$ denote the collection of closures of the balls of \mathcal{G}_i . Then let \mathcal{F}_m be those balls of \mathcal{F} , such that if B is one of these balls, \overline{B} has empty intersection with every closed ball of $\overline{\mathcal{G}_i}$ for each $i \leq m-1$. Then using Lemma 4.5.2, let \mathcal{G}_m be a maximal collection of balls from \mathcal{F}_m with the property that each ball has radius larger than $(\frac{2}{3})^m M$ and their closures are disjoint. Let $\mathcal{G} \equiv \cup_{k=1}^{\infty} \mathcal{G}_k$. Thus the closures of balls in \mathcal{G} are disjoint. Let $\mathbf{x} \in B(\mathbf{p}, r) \in \mathcal{F} \setminus \mathcal{G}$. Choose m such that

$$\left(\frac{2}{3}\right)^m M < r \leq \left(\frac{2}{3}\right)^{m-1} M$$

Then $\overline{B(\mathbf{p}, r)}$ must have nonempty intersection with the closure of some ball from $\mathcal{G}_1 \cup \dots \cup \mathcal{G}_m$ because if it didn't, then \mathcal{G}_m would fail to be maximal. Denote by $B(\mathbf{p}_0, r_0)$ a ball in $\mathcal{G}_1 \cup \dots \cup \mathcal{G}_m$ whose closure has nonempty intersection with $\overline{B(\mathbf{p}, r)}$. Thus both

$$r_0, r > \left(\frac{2}{3}\right)^m M, \text{ so } r \leq \left(\frac{2}{3}\right)^{m-1} M < \frac{3}{2} r_0$$

Consider the picture, in which $\mathbf{w} \in \overline{B(\mathbf{p}_0, r_0)} \cap \overline{B(\mathbf{p}, r)}$.



Then for $x \in \overline{B(p, r)}$,

$$\begin{aligned} \|x - p_0\| &\leq \|x - p\| + \|p - w\| + \overbrace{\|w - p_0\|}^{\leq r_0} \\ &\leq r + r + r_0 \leq 2 \overbrace{\left(\frac{2}{3}\right)^{m-1} M}^{< \frac{3}{2}r_0} + r_0 \leq 2 \left(\frac{3}{2}r_0\right) + r_0 \leq 4r_0 \end{aligned}$$

Thus $B(p, r)$ is contained in $\overline{B(p_0, 4r_0)}$. It follows that the closures of the balls of \mathcal{G} are disjoint and the set $\{\hat{B} : B \in \mathcal{G}\}$ covers A . ■

Note that this theorem does not depend on the underlying space being finite dimensional. However, it is typically used in this setting. The next theorem of Besicovitch depends in an essential way on X being finite dimensional because it exploits compactness and various constants originate explicitly from this compactness. However, no effort is being made here to give the most general conditions under which such covering theorems hold.

4.5.2 Besicovitch Covering Theorem

The covering theorems will have applications to measure theory presented later. In contrast to the Vitali covering theorem, one does not enlarge the balls in the Besicovitch covering theorem. This is extremely useful in the notion of general differentiation theorems for measures other than Lebesgue measure. The proof of this major result has to do with counting the number of times various balls can intersect. These estimates are used along with the pigeon hole principle to prove the result. This principle says that if you have n holes and $m > n$ pigeons, each of which must go in a hole, then some hole has more than one pigeon. In what follows x will continue to be in a normed linear space $(X, \|\cdot\|)$ of dimension p . This covering theorem is one of the most amazing and insightful ideas that I have ever encountered. It is simultaneously elegant, elementary and profound. This section is an attempt to present this wonderful result.

Here is a sequence of balls from \mathcal{F} in the case that the set of centers of these balls is bounded. I will denote by $r(B_k)$ the radius of a ball B_k .

A construction of a sequence of balls

Lemma 4.5.4 *Let \mathcal{F} be a nonempty set of nonempty balls in X with*

$$\sup \{\text{diam}(B) : B \in \mathcal{F}\} = D < \infty$$

and let A denote the set of centers of these balls. Suppose A is bounded. Define a sequence of balls from \mathcal{F} , $\{B_j\}_{j=1}^J$ where $J \leq \infty$ such that

$$r(B_1) > \frac{3}{4} \sup \{r(B) : B \in \mathcal{F}\} \quad (4.26)$$

and if

$$A_m \equiv A \setminus (\cup_{i=1}^m B_i) \neq \emptyset, \quad (4.27)$$

then $B_{m+1} \in \mathcal{F}$ is chosen with center in A_m such that

$$r(B_m) > r(B_{m+1}) > \frac{3}{4} \sup \{r : B(a, r) \in \mathcal{F}, a \in A_m\}. \quad (4.28)$$

Then letting $B_j = B(a_j, r_j)$, this sequence satisfies $\{B(a_j, r_j/3)\}_{j=1}^J$ are disjoint.

$$A \subseteq \cup_{i=1}^J B_i. \quad (4.29)$$

Proof: First note that B_{m+1} can be chosen as in 4.28. This is because the A_m are decreasing and so

$$\begin{aligned} & \frac{3}{4} \sup \{r : B(a, r) \in \mathcal{F}, a \in A_m\} \\ & \leq \frac{3}{4} \sup \{r : B(a, r) \in \mathcal{F}, a \in A_{m-1}\} < r(B_m) \end{aligned}$$

Thus the $r(B_k)$ are strictly decreasing and so no B_k contains a center of any other B_j .

If $x \in B(a_j, r_j/3) \cap B(a_i, r_i/3)$ where these balls are two which are chosen by the above scheme such that $j > i$, then from what was just shown

$$\|a_j - a_i\| \leq \|a_j - x\| + \|x - a_i\| \leq \frac{r_j}{3} + \frac{r_i}{3} \leq \left(\frac{1}{3} + \frac{1}{3}\right) r_i = \frac{2}{3} r_i < r_i$$

and this contradicts the construction because a_j is not covered by $B(a_i, r_i)$.

Finally consider the claim that $A \subseteq \cup_{i=1}^J B_i$. Pick B_1 satisfying 4.26. If

$$B_1, \dots, B_m$$

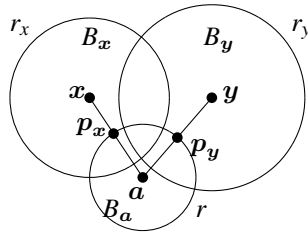
have been chosen, and A_m is given in 4.27, then if $A_m = \emptyset$, it follows $A \subseteq \cup_{i=1}^m B_i$. Set $J = m$.

Now let a be the center of $B_a \in \mathcal{F}$. If $a \in A_m$ for all m , (That is a does not get covered by the B_i .) then $r_{m+1} \geq \frac{3}{4} r(B_a)$ for all m , a contradiction since the balls $B(a_j, \frac{r_j}{3})$ are disjoint and A is bounded, implying that $r_j \rightarrow 0$. Thus a must fail to be in some A_m which means it was covered by some ball in the sequence. ■

The covering theorem is obtained by estimating how many B_j can intersect B_k for $j < k$. The thing to notice is that from the construction, no B_j contains the center of another B_i . Also, the $r(B_k)$ is a decreasing sequence.

Let $\alpha > 1$. There are two cases for an intersection. Either $r(B_j) \geq \alpha r(B_k)$ or $\alpha r(B_k) > r(B_j) > r(B_k)$.

First consider the case where we have a ball $B(a, r)$ intersected with other balls of radius larger than αr such that none of the balls contains the center of any other. This is illustrated in the following picture with two balls. This has to do with estimating the number of B_j for $j \leq k$ where $r(B_j) \geq \alpha r(B_k)$.



Imagine projecting the center of each big ball as in the above picture onto the surface of the given ball, assuming the given ball has radius 1. By scaling the balls, you could reduce to this case that the given ball has radius 1. Then from geometric reasoning, there should be a lower bound to the distance between these two projections depending on dimension. Thus there is an estimate on how many large balls can intersect the given ball with no ball containing a center of another one.

Intersections with relatively big balls

Lemma 4.5.5 *Let the balls B_a, B_x, B_y be as shown, having radii r, r_x, r_y respectively. Suppose the centers of B_x and B_y are not both in any of the balls shown, and suppose $r_y \geq r_x \geq \alpha r$ where α is a number larger than 1. Also let $P_x \equiv a + r \frac{x-a}{\|x-a\|}$ with P_y being defined similarly. Then it follows that $\|P_x - P_y\| \geq \frac{\alpha-1}{\alpha+1}r$. There exists a constant $L(p, \alpha)$ depending on α and the dimension, such that if B_1, \dots, B_m are all balls such that any pair are in the same situation relative to B_a as B_x and B_y , then $m \leq L(p, \alpha)$.*

Proof: From the definition,

$$\begin{aligned}
 \|P_x - P_y\| &= r \left\| \frac{x-a}{\|x-a\|} - \frac{y-a}{\|y-a\|} \right\| \\
 &= r \left\| \frac{(x-a)\|y-a\| - (y-a)\|x-a\|}{\|x-a\|\|y-a\|} \right\| \\
 &= r \left\| \frac{\|y-a\|(x-y) + (y-a)(\|y-a\| - \|x-a\|)}{\|x-a\|\|y-a\|} \right\| \\
 &\geq r \frac{\|x-y\|}{\|x-a\|} - r \frac{\|y-a\| |\|y-a\| - \|x-a\||}{\|x-a\|\|y-a\|} \\
 &= r \frac{\|x-y\|}{\|x-a\|} - \frac{r}{\|x-a\|} \left| \|y-a\| - \|x-a\| \right|. \tag{4.30}
 \end{aligned}$$

There are two cases. First suppose that $\|y-a\| - \|x-a\| \geq 0$. Then the above

$$= r \frac{\|x-y\|}{\|x-a\|} - \frac{r}{\|x-a\|} \|y-a\| + r.$$

From the assumptions, $\|x-y\| \geq r_y$ and also $\|y-a\| \leq r + r_y$. Hence the above

$$\begin{aligned}
 &\geq r \frac{r_y}{\|x-a\|} - \frac{r}{\|x-a\|} (r + r_y) + r = r - r \frac{r}{\|x-a\|} \\
 &\geq r \left(1 - \frac{r}{\|x-a\|} \right) \geq r \left(1 - \frac{r}{r_x} \right) \geq r \left(1 - \frac{1}{\alpha} \right) \geq r \frac{\alpha-1}{\alpha+1}.
 \end{aligned}$$

The other case is that $\|y-a\| - \|x-a\| < 0$ in 4.30. Then in this case 4.30 equals

$$\begin{aligned}
 &= r \left(\frac{\|x-y\|}{\|x-a\|} - \frac{1}{\|x-a\|} (\|x-a\| - \|y-a\|) \right) \\
 &= \frac{r}{\|x-a\|} (\|x-y\| - (\|x-a\| - \|y-a\|))
 \end{aligned}$$

Then since $\|\mathbf{x} - \mathbf{a}\| \leq r + r_x$, $\|\mathbf{x} - \mathbf{y}\| \geq r_y$, $\|\mathbf{y} - \mathbf{a}\| \geq r_y$, and remembering that $r_y \geq r_x \geq \alpha r$,

$$\begin{aligned} &\geq \frac{r}{r_x + r} (r_y - (r + r_x) + r_y) \geq \frac{r}{r_x + r} (r_y - (r + r_y) + r_y) \\ &\geq \frac{r}{r_x + r} (r_y - r) \geq \frac{r}{r_x + r} (r_x - r) \geq \frac{r}{r_x + \frac{1}{\alpha} r_x} \left(r_x - \frac{1}{\alpha} r_x \right) \\ &= \frac{r}{1 + (1/\alpha)} (1 - 1/\alpha) = \frac{\alpha - 1}{\alpha + 1} r \end{aligned}$$

Replacing r with something larger, $\frac{1}{\alpha} r_x$ is justified by the observation that $x \rightarrow \frac{\alpha - x}{\alpha + x}$ is decreasing. This proves the estimate between P_x and P_y .

Finally, in the case of the balls B_i having centers at \mathbf{x}_i , then as above, let $P_{\mathbf{x}_i} = \mathbf{a} + r \frac{\mathbf{x}_i - \mathbf{a}}{\|\mathbf{x}_i - \mathbf{a}\|}$. Then $(P_{\mathbf{x}_i} - \mathbf{a}) r^{-1}$ is on the unit sphere having center $\mathbf{0}$. Furthermore,

$$\|(P_{\mathbf{x}_i} - \mathbf{a}) r^{-1} - (P_{\mathbf{y}_i} - \mathbf{a}) r^{-1}\| = r^{-1} \|P_{\mathbf{x}_i} - P_{\mathbf{y}_i}\| \geq r^{-1} r \frac{\alpha - 1}{\alpha + 1} = \frac{\alpha - 1}{\alpha + 1}.$$

How many points on the unit sphere can be pairwise this far apart? The unit sphere is compact and so there exists a $\frac{1}{4} \left(\frac{\alpha - 1}{\alpha + 1} \right)$ net having $L(p, \alpha)$ points. Thus m cannot be any larger than $L(p, \alpha)$ because if it were, then by the pigeon hole principal, two of the points $(P_{\mathbf{x}_i} - \mathbf{a}) r^{-1}$ would lie in a single ball $B(p, \frac{1}{4} \left(\frac{\alpha - 1}{\alpha + 1} \right))$ so they could not be $\frac{\alpha - 1}{\alpha + 1}$ apart. ■

The above lemma has to do with balls which are relatively large intersecting a given ball. Next is a lemma which has to do with relatively small balls intersecting a given ball. First is another lemma.

Lemma 4.5.6 *Let $\Gamma > 1$ and $B(\mathbf{a}, \Gamma r)$ be a ball and suppose $\{B(\mathbf{x}_i, r_i)\}_{i=1}^m$ are balls contained in $B(\mathbf{a}, \Gamma r)$ such that $r \leq r_i$ and none of these balls contains the center of another ball. Then there is a constant $M(p, \Gamma)$ such that $m \leq M(p, \Gamma)$.*

Proof: Let $\mathbf{z}_i = \mathbf{x}_i - \mathbf{a}$. Then $B(\mathbf{z}_i, r_i)$ are balls contained in $B(\mathbf{0}, \Gamma r)$ with no ball containing a center of another. Then $B(\frac{\mathbf{z}_i}{\Gamma r}, \frac{r_i}{\Gamma r})$ are balls in $B(\mathbf{0}, 1)$ with no ball containing the center of another. By compactness, there is a $\frac{1}{8\Gamma}$ net for $\overline{B(\mathbf{0}, 1)}$, $\{\mathbf{y}_i\}_{i=1}^{M(p, \Gamma)}$. Thus the balls $B(\mathbf{y}_i, \frac{1}{8\Gamma})$ cover $\overline{B(\mathbf{0}, 1)}$. If $m \geq M(p, \Gamma)$, then by the pigeon hole principle, one of these $B(\mathbf{y}_i, \frac{1}{8\Gamma})$ would contain some $\frac{\mathbf{z}_i}{\Gamma r}$ and $\frac{\mathbf{z}_j}{\Gamma r}$ which requires $\|\frac{\mathbf{z}_i}{\Gamma r} - \frac{\mathbf{z}_j}{\Gamma r}\| \leq \frac{1}{4\Gamma} < \frac{r_j}{4\Gamma r}$ so $\frac{\mathbf{z}_i}{\Gamma r} \in B(\frac{\mathbf{z}_j}{\Gamma r}, \frac{r_j}{\Gamma r})$. Thus $m \leq M(p, \Gamma, \Gamma)$. ■

Intersections with small balls

Lemma 4.5.7 *Let B be a ball having radius r and suppose B has nonempty intersection with the balls B_1, \dots, B_m having radii r_1, \dots, r_m respectively, and as before, no B_i contains the center of any other and the centers of the B_i are not contained in B . Suppose $\alpha > 1$ and $r \leq \min(r_1, \dots, r_m)$, each $r_i < \alpha r$. Then there exists a constant $M(p, \alpha)$ such that $m \leq M(p, \alpha)$.*

Proof: Let $B = B(\mathbf{a}, r)$. Then each B_i is contained in $B(\mathbf{a}, 2r + \alpha r + \alpha r)$. This is because if $\mathbf{y} \in B_i \equiv B(\mathbf{x}_i, r_i)$,

$$\|\mathbf{y} - \mathbf{a}\| \leq \|\mathbf{y} - \mathbf{x}_i\| + \|\mathbf{x}_i - \mathbf{a}\| \leq r_i + r + r_i < 2r + \alpha r + \alpha r$$

Thus B_i does not contain the center of any other B_j , these balls are each contained in $B(a, r(2\alpha + 2))$, and each radius is at least as large as r . By Lemma 4.5.6 there is a constant $M(p, \alpha)$ such that $m \leq M(p, \alpha)$. ■

Now here is the Besicovitch covering theorem. In the proof, we are considering the sequence of balls described above.

Theorem 4.5.8 *There exists a constant N_p , depending only on p with the following property. If \mathcal{F} is any collection of nonempty balls in X with*

$$\sup \{ \text{diam}(B) : B \in \mathcal{F} \} < D < \infty$$

and if A is the set of centers of the balls in \mathcal{F} , then there exist subsets of \mathcal{F} , $\mathcal{H}_1, \dots, \mathcal{H}_{N_p}$, such that each \mathcal{H}_i is a countable collection of disjoint balls from \mathcal{F} (possibly empty) and

$$A \subseteq \bigcup_{i=1}^{N_p} \mathcal{H}_i \cup \{B : B \in \mathcal{F}\}.$$

Proof: To begin with, suppose A is bounded. Let $L(p, \alpha)$ be the constant of Lemma 4.5.5 and let $M_p = L(p, \alpha) + M(p, \alpha) + 1$. Define the following sequence of subsets of \mathcal{F} , $\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_{M_p}$. Referring to the sequence $\{B_k\}$ considered in Lemma 4.5.4, let $B_1 \in \mathcal{G}_1$ and if B_1, \dots, B_m have been assigned, each to a \mathcal{G}_i , place B_{m+1} in the first \mathcal{G}_j such that B_{m+1} intersects no set already in \mathcal{G}_j . The existence of such a j follows from Lemmas 4.5.5 and 4.5.7 and the pigeon hole principle. Here is why. B_{m+1} can intersect at most $L(p, \alpha)$ sets of $\{B_1, \dots, B_m\}$ which have radii at least as large as $\alpha r(B_{m+1})$ thanks to Lemma 4.5.5. It can intersect at most $M(p, \alpha)$ sets of $\{B_1, \dots, B_m\}$ which have radius smaller than $\alpha r(B_{m+1})$ thanks to Lemma 4.5.7. Thus each \mathcal{G}_j consists of disjoint sets of \mathcal{F} and the set of centers is covered by the union of these \mathcal{G}_j . This proves the theorem in case the set of centers is bounded.

Now let $R_1 = B(0, 5D)$ and if R_m has been chosen, let

$$R_{m+1} = B(0, (m+1)5D) \setminus R_m$$

Thus, if $|k - m| \geq 2$, no ball from \mathcal{F} having nonempty intersection with R_m can intersect any ball from \mathcal{F} which has nonempty intersection with R_k . This is because all these balls have radius less than D . Now let $A_m \equiv A \cap R_m$ and apply the above result for a bounded set of centers to those balls of \mathcal{F} which intersect R_m to obtain sets of disjoint balls $\mathcal{G}_1(R_m), \mathcal{G}_2(R_m), \dots, \mathcal{G}_{M_p}(R_m)$ covering A_m . Then simply define $\mathcal{G}'_j \equiv \bigcup_{k=1}^{\infty} \mathcal{G}_j(R_{2k})$, $\mathcal{G}_j \equiv \bigcup_{k=1}^{\infty} \mathcal{G}_j(R_{2k-1})$. Let $N_p = 2M_p$ and

$$\{\mathcal{H}_1, \dots, \mathcal{H}_{N_p}\} \equiv \{\mathcal{G}'_1, \dots, \mathcal{G}'_{M_p}, \mathcal{G}_1, \dots, \mathcal{G}_{M_p}\}$$

Note that the balls in \mathcal{G}'_j are disjoint. This is because those in $\mathcal{G}_j(R_{2k})$ are disjoint and if you consider any ball in $\mathcal{G}_j(R_{2m})$, it cannot intersect a ball of $\mathcal{G}_j(R_{2k})$ for $m \neq k$ because $|2k - 2m| \geq 2$. Similar considerations apply to the balls of \mathcal{G}_j . ■

Of course, you could pick a particular α . If you make α larger, $L(p, \alpha)$ should get smaller and $M(p, \alpha)$ should get larger. Obviously one could explore this at length to try and get a best choice of α .

4.6 Exercises

1. Let V be a vector space with basis $\{v_1, \dots, v_n\}$. For $v \in V$, denote its coordinate vector as $\mathbf{v} = (\alpha_1, \dots, \alpha_n)$ where $v = \sum_{k=1}^n \alpha_k v_k$. Now define

$$\|v\| \equiv \max \{|\alpha_k| : k = 1, \dots, n\}.$$

Show that this is a norm on V .

2. Let $(X, \|\cdot\|)$ be a normed linear space. You can let it be $(\mathbb{R}^n, |\cdot|)$ if you like. Recall $|x|$ is the usual magnitude of a vector given by $|x| = \sqrt{\sum_{k=1}^n |x_k|^2}$. A set A is said to be **convex** if whenever $\mathbf{x}, \mathbf{y} \in A$ the line segment determined by these points given by $t\mathbf{x} + (1-t)\mathbf{y}$ for $t \in [0, 1]$ is also in A . Show that every open or closed ball is convex. Remember a closed ball is $D(\mathbf{x}, r) \equiv \{\hat{\mathbf{x}} : \|\hat{\mathbf{x}} - \mathbf{x}\| \leq r\}$ while the open ball is $B(\mathbf{x}, r) \equiv \{\hat{\mathbf{x}} : \|\hat{\mathbf{x}} - \mathbf{x}\| < r\}$. This should work just as easily in any normed linear space with any norm.

3. This problem is for those who have had a course in Linear algebra. A vector \mathbf{v} is in the convex hull of S if there are finitely many vectors of S , $\{v_1, \dots, v_m\}$ and nonnegative scalars $\{t_1, \dots, t_m\}$ such that $\mathbf{v} = \sum_{k=1}^m t_k v_k$, $\sum_{k=1}^m t_k = 1$. Such a linear combination is called a convex combination. Suppose now that $S \subseteq V$, a vector space of dimension n . Show that if $\mathbf{v} = \sum_{k=1}^m t_k v_k$ is a vector in the convex hull for $m > n + 1$, then there exist other nonnegative scalars $\{t'_k\}$ summing to 1 such that $\mathbf{v} = \sum_{k=1}^{m-1} t'_k v_k$. Thus every vector in the convex hull of S can be obtained as a convex combination of at most $n + 1$ points of S . This incredible result is in Rudin [51]. Convexity is more a geometric property than a topological property. **Hint:** Consider $L: \mathbb{R}^m \rightarrow V \times \mathbb{R}$ defined by $L(\mathbf{a}) \equiv (\sum_{k=1}^m a_k v_k, \sum_{k=1}^m a_k)$. Explain why $\ker(L) \neq \{0\}$. This will involve observing that \mathbb{R}^m has higher dimension than $V \times \mathbb{R}$. Thus L cannot be one to one because one to one functions take linearly independent sets to linearly independent sets and you can't have a linearly independent set with more than $n + 1$ vectors in $V \times \mathbb{R}$. Next, letting $\mathbf{a} \in \ker(L) \setminus \{0\}$ and $\lambda \in \mathbb{R}$, note that $\lambda \mathbf{a} \in \ker(L)$. Thus for all $\lambda \in \mathbb{R}$, $\mathbf{v} = \sum_{k=1}^m (t_k + \lambda a_k) v_k$. Now vary λ till some $t_k + \lambda a_k = 0$ for some $a_k \neq 0$. You can assume each $t_k > 0$ since otherwise, there is nothing to show. This is a really nice result because it can be used to show that the convex hull of a compact set is also compact. You might try to show this if you feel like it.

4. Show that the usual norm in \mathbb{F}^n given by $|\mathbf{x}| = (\mathbf{x}, \mathbf{x})^{1/2}$ satisfies the following identities, the first of them being the parallelogram identity and the second being the polarization identity.

$$\begin{aligned} |\mathbf{x} + \mathbf{y}|^2 + |\mathbf{x} - \mathbf{y}|^2 &= 2|\mathbf{x}|^2 + 2|\mathbf{y}|^2 \\ \operatorname{Re}(\mathbf{x}, \mathbf{y}) &= \frac{1}{4} (|\mathbf{x} + \mathbf{y}|^2 - |\mathbf{x} - \mathbf{y}|^2) \end{aligned}$$

Show that these identities hold in any inner product space, not just \mathbb{F}^n .

5. Suppose K is a compact subset of (X, d) a metric space. Also let \mathcal{C} be an open cover of K . Show that there exists $\delta > 0$ such that for all $x \in K$, $B(x, \delta)$ is contained in a single set of \mathcal{C} . This number is called a Lebesgue number. **Hint:** For each $x \in K$, there exists $B(x, \delta_x)$ such that this ball is contained in a set of \mathcal{C} . Now consider

the balls $\left\{B\left(x, \frac{\delta_x}{2}\right)\right\}_{x \in K}$. Finitely many of these cover K . $\left\{B\left(x_i, \frac{\delta_{x_i}}{2}\right)\right\}_{i=1}^n$ Now consider what happens if you let $\delta \leq \min\left\{\frac{\delta_{x_i}}{2}, i = 1, 2, \dots, n\right\}$. Explain why this works. You might draw a picture to help get the idea.

6. Suppose \mathcal{C} is a set of compact sets in a metric space (X, d) and suppose that the intersection of **every** finite subset of \mathcal{C} is nonempty. This is called the **finite intersection property**. Show that $\cap \mathcal{C}$, the intersection of all sets of \mathcal{C} is nonempty. This particular result is enormously important. **Hint:** You could let \mathcal{U} denote the set $\{K^C : K \in \mathcal{C}\}$. If $\cap \mathcal{C}$ is empty, then its complement is $\cup \mathcal{U} = X$. Picking $K \in \mathcal{C}$, it follows that \mathcal{U} is an open cover of K . $K \subseteq \cup_{i=1}^m K_i^C = \left(\cap_{i=1}^m K_i\right)^C$ Therefore, you would need to have $\{K_1^C, \dots, K_m^C\}$ is a cover of K . In other words, Now what does this say about the intersection of K with these K_i ?
7. If (X, d) is a compact metric space and $f : X \rightarrow Y$ is continuous where (Y, ρ) is another metric space, show that if f is continuous on X , then it is uniformly continuous. Recall that this means that if $\varepsilon > 0$ is given, then there exists $\delta > 0$ such that if $d(x, \hat{x}) < \delta$, then $\rho(f(x), f(\hat{x})) < \varepsilon$. Compare with the definition of continuity. **Hint:** If this is not so, then there exists $\varepsilon > 0$ and x_n, \hat{x}_n such that $d(x_n, \hat{x}_n) < 1/n$ but $\rho(f(x_n), f(\hat{x}_n)) \geq \varepsilon$. Now use compactness to get a contradiction.
8. Prove the above problem using another approach. Use the existence of the Lebesgue number in Problem 5 to prove continuity on a compact set K implies uniform continuity on this set. **Hint:** Consider $\mathcal{C} \equiv \{f^{-1}(B(f(x), \varepsilon/2)) : x \in X\}$. This is an open cover of X . Let δ be a Lebesgue number for this open cover. Suppose $d(x, \hat{x}) < \delta$. Then both x, \hat{x} are in $B(x, \delta)$ and so both are in $f^{-1}(B(f(\bar{x}), \frac{\varepsilon}{2}))$. Hence

$$\rho(f(x), f(\bar{x})) < \frac{\varepsilon}{2}, \rho(f(\hat{x}), f(\bar{x})) < \frac{\varepsilon}{2}.$$

Now consider the triangle inequality.

9. Let X be a vector space. A Hamel basis is a subset of X, Λ such that every vector of X can be written as a finite linear combination of vectors of Λ and the vectors of Λ are linearly independent in the sense that if $\{x_1, \dots, x_n\} \subseteq \Lambda$ and $\sum_{k=1}^n c_k x_k = 0$ then each $c_k = 0$. Using the Hausdorff maximal theorem, show that every non-zero vector space has a Hamel basis. **Hint:** Let $x_1 \neq 0$. Let \mathcal{F} denote the collection of subsets of X, Λ containing x_1 with the property that the vectors of Λ are linearly independent. Partially order \mathcal{F} by set inclusion and consider the union of a maximal chain.
10. Suppose X is a nonzero real or complex normed linear space and let

$$V = \text{span}(w_1, \dots, w_m)$$

where $\{w_1, \dots, w_m\}$ is a linearly independent set of vectors of X . Show that V is a closed subspace of X with $V \subsetneq X$. First explain why Theorem 4.2.11 implies any finite dimensional subspace of X can be written this way. **Hint:** You might want to use something like Lemma 4.4.7 to show this.

11. Suppose X is a normed linear space and its dimension is either infinite or greater than m where $V \equiv \text{span}(w_1, \dots, w_m)$ for $\{w_1, \dots, w_m\}$ an independent set of vectors of X .

Show $X \setminus V$ is a dense open subset of X which is equivalent to V containing no ball $B(v, r)$, $\{w : \|w - v\| < r\}$. **Hint:** If $B(x, r)$ is contained in V , then show, that since V is a subspace, $B(0, r)$ is contained in V . Then show this implies $X \subseteq V$ which is not the case.

12. Show that if (X, d) is a metric space and H, K are disjoint closed sets, there are open sets U_H, U_K such that $H \subseteq U_H, K \subseteq U_K$ and $U_H \cap U_K = \emptyset$. **Hint:** Let $k \in K$. Explain why $\text{dist}(k, H) \equiv \inf\{\|k - h\| : h \in H\} \equiv 2\delta_k > 0$. Now consider $U_K \equiv \cup_{k \in K} B(k, \delta_k)$. Do something similar for $h \in H$ and consider $U_H \equiv \cup_{h \in H} B(h, \delta_h)$.
13. If, in a metric space, $B(p, \delta)$ is a ball, show that

$$\overline{B(p, \delta)} \subseteq D(p, \delta) \equiv \{x : \|x - p\| \leq \delta\}$$

Now suppose (X, d) is a complete metric space and $U_n, n \in \mathbb{N}$ is a dense open set in X . Also let W be any nonempty open set. Show there exists a ball $B_1 \equiv B(p_1, r_1)$ having radius smaller than 2^{-1} such that $\overline{B_1} \subseteq U_1 \cap W$. Next show there exists $B_2 \equiv B(p_2, r_2)$ such that $\overline{B_2} \subseteq B_1 \cap U_2 \cap W$ with the radius of B_2 less than 2^{-2} . Continue this way. Explain why $\{p_n\}_{n=1}^\infty$ is a Cauchy sequence converging to some $p \in W \cap (\cup_{n=1}^\infty U_n)$. This is the very important Baire theorem which says that in a complete metric space, the intersection of dense open sets is dense.

14. Suppose you have a complete normed linear space, $(X, \|\cdot\|)$. Use the above problems leading to the Baire theorem in 13 to show that if \mathcal{B} is a Hamel basis for X , then \mathcal{B} cannot be countable. **Hint:** If $\mathcal{B} = \{v_i\}_{i=1}^\infty$, consider $V_n \equiv \text{span}(v_1, \dots, v_n)$. Then use a problem listed above to argue that V_n^C is a dense open set. Now apply Problem 13. This shows why the idea of a Hamel basis often fails to be very useful whereas, in finite dimensional settings, it is just what is needed.
15. In any complete normed linear space which is infinite dimensional, show the unit ball is not compact. Do this by showing the existence of a sequence which cannot have a convergent subsequence. **Hint:** Pick $\|x_1\| = 1$. Suppose x_1, \dots, x_n have been chosen, each $\|x_k\| = 1$. Then there is $x \notin \text{span}(x_1, \dots, x_n) \equiv V_n$. Now consider v such that $\|x - v\| \leq \frac{3}{2} \text{dist}(x, V_n)$. Then argue that for $k \leq n$,

$$\left\| \frac{x - v}{\|x - v\|} - x_k \right\| = \left\| \frac{x - \left(v + \overbrace{\|x - v\| x_k}^{\in V_n} \right)}{\|x - v\|} \right\| \geq \frac{\text{dist}(x, V_n)}{(3/2) \text{dist}(x, V_n)} = \frac{2}{3}$$

16. Let X be a complete inner product space. Let \mathcal{F} denote subsets $\beta \subseteq X$ such that whenever $x, y \in \beta, (x, y) = 0$ if $x \neq y$ and $(x, x) = 1$ if $x = y$. Thus these β are orthonormal sets. Show there exists a maximal orthonormal set. If X is separable, show that this maximal orthonormal set is countable. **Hint:** Use the Hausdorff maximal theorem. The next few problems involve linear algebra.
17. Let X be a real inner product space and let $\{v_1, \dots, v_n\}$ be vectors in X . Let G be the $n \times n$ matrix $G_{ij} \equiv (v_i, v_j)$. Show that G^{-1} exists if and only if $\{v_1, \dots, v_n\}$ is linearly independent. G is called the Grammian or the metric tensor.

18. \uparrow Let X be as above, a real inner product space, and let $V \equiv \text{span}(v_1, \dots, v_n)$. Let $u \in X$ and $z \in V$. Show that $|u - z| = \inf\{|u - v| : v \in V\}$ if and only if $(u - z, v_i) = 0$ for all v_i . Note that the v_i might not be linearly independent. Also show that $|u - z|^2 = |u|^2 - (z, u)$.
19. \uparrow Let G be the matrix of Problem 17 where $\{v_1, \dots, v_n\}$ is linearly independent and $V \equiv \text{span}(v_1, \dots, v_n) \subseteq X$, an inner product space. Let $x \equiv \sum_i x^i v_i, y \equiv \sum_i y^i v_i$ be two vectors of V . Show that $(x, y) = \sum_{i,j} x^i G_{ij} y^j$. Show that $z \equiv \sum_i z^i v_i, z$ is closest to $u \in X$ if and only if for all $i = 1, \dots, n, (u, v_i) = \sum_j G_{ij} z^j$. This gives a system of linear equations which must be satisfied by the z^i in order that z just given is the best approximation to u . Next show that there exists such a solution thanks to Problem 17 which says that the matrix G is invertible, and if G^{-1} has i, j^{th} component G^{ij} , one finds that $\sum_j G^{ij} (u, v_j) = z^i$.
20. \uparrow In the situation of the above problems, suppose A is an $m \times n$ matrix. Use Problem 18 to show that for $y \in \mathbb{R}^m$, there always exists a solution x to the system of equations $A^T y = A^T A x$. Explain how this is in a sense the best you can do to solve $y = A x$ even though this last system of equations might not have a solution. Here A^T is the transpose of the matrix A . The equations $A^T y = A^T A x$ are called the normal equations for the least squares problem. **Hint:** Verify that $(A^T y, x) = (y, A x)$. Let the subspace V be $A(\mathbb{R}^n)$, the vectors spanning it being $\{A e_1, \dots, A e_n\}$. From the above problem, there exists $A x$ in V which is closest to y . Now use the characterization of this vector $(y - A x, A z) = 0$ for all $z \in \mathbb{R}^n, A z$ being a generic vector in $A(\mathbb{R}^n)$.
21. \uparrow As an example of an inner product space, consider $C([0, 1])$ with the inner product $\int_0^1 f(x) g(x) dx$ where this is the ordinary integral from calculus. Abusing notation, let $\{x^{p_1}, \dots, x^{p_n}\}$ with $-\frac{1}{2} < p_1 < \dots < p_n$ be functions, (vectors) in $C([0, 1])$. Verify that these vectors are linearly independent. **Hint:** You might want to use the Cauchy identity, Theorem 1.9.28.
22. \uparrow As above, if $\{v_1, \dots, v_n\}$ is linearly independent, the Gramian is $G = G(v_1, \dots, v_n)$, $G_{ij} \equiv (v_i, v_j)$, then if $u \notin \text{span}(v_1, \dots, v_n) \equiv V$ you could consider $G(v_1, \dots, v_n, u)$. Then if $d \equiv \min\{|u - v| : v \in \text{span}(v_1, \dots, v_n)\}$, show that $d^2 = \frac{\det G(v_1, \dots, v_n, u)}{\det G(v_1, \dots, v_n)}$. Justify the following steps. Letting z be the closest point of V to u , from the above, $(u - \sum_{i=1}^n z^i v_i, v_p) = 0$ for each v_p and so

$$(u, v_p) = \sum_{i=1}^n (v_p, v_i) z^i \quad (*)$$

Also, since $(u - z, v) = 0$ for all $v \in V, |u|^2 = |u - z + z|^2 = |u - z|^2 + |z|^2$ so

$$\begin{aligned} |u|^2 &= \left| u - \sum_{i=1}^n z^i v_i \right|^2 + \left| \sum_{i=1}^n z^i v_i \right|^2 = d^2 + \left| \sum_{i=1}^n z^i v_i \right|^2 \\ &= d^2 + \sum_j \overbrace{\sum_i (v_j, v_i) z^i z^j}^{=(u, v_j)} = d^2 + \sum_j (u, v_j) z^j \end{aligned}$$

$$= d^2 + \mathbf{y}^T \mathbf{z}, \mathbf{y} \equiv ((u, v_1), \dots, (u, v_n))^T, \mathbf{z} \equiv (z^1, \dots, z^n)^T$$

From *, $G\mathbf{z} = \mathbf{y}$, $\begin{pmatrix} G(v_1, \dots, v_n) & \mathbf{0} \\ \mathbf{y}^T & 1 \end{pmatrix} \begin{pmatrix} \mathbf{z} \\ d^2 \end{pmatrix} = \begin{pmatrix} \mathbf{y} \\ \|u\|^2 \end{pmatrix}$. Now use Cramer's rule to solve for d^2 and get

$$d^2 = \frac{\det \begin{pmatrix} G(v_1, \dots, v_n) & \mathbf{y} \\ \mathbf{y}^T & \|u\|^2 \end{pmatrix}}{\det(G(v_1, \dots, v_n))} \equiv \frac{\det G(v_1, \dots, v_n, u)}{\det G(v_1, \dots, v_n)}$$

23. In the situation of Problem 21, let $f_k(x) \equiv x^k$ and let $V \equiv \text{span}(f_{p_1}, \dots, f_{p_n})$. give an estimate for the distance d between f_m and V for m a nonnegative integer and as in the above problem $-\frac{1}{2} < p_1 < \dots < p_n$. Use Theorem 1.9.28 in the appendix and the above problem with $v_i \equiv f_{p_i}$ and $v_{n+1} \equiv f_m$. Justify the following manipulations.

The numerator in the above formula for the distance is of the form $\frac{\prod_{j < i \leq n+1} (p_i - p_j)^2}{\prod_{i, j \leq n+1} (p_i + p_j + 1)}$

$$= \frac{\prod_{j < i \leq n} (p_i - p_j)^2 \prod_{j \leq n} (m - p_j)^2}{\prod_{i, j \leq n} (p_i + p_j + 1) \prod_{i=1}^n (p_i + m + 1) \prod_{j=1}^n (p_j + m + 1) (2m + 1)}$$

While $G(f_{p_1}, \dots, f_{p_n}) = \frac{\prod_{j < i \leq n} (p_i - p_j)^2}{\prod_{i, j \leq n} (p_i + p_j + 1)}$. Thus $d = \frac{\prod_{j \leq n} |m - p_j|}{\prod_{i=1}^n (p_i + m + 1) (\sqrt{2m + 1})}$.

24. Suppose $\sum_{k=0}^n a_k t^k = 0$ for each $t \in (-\delta, \delta)$ where $a_k \in X$, a linear space. Show that each $a_k = 0$.
25. Suppose $A \subseteq \mathbb{R}^p$ is covered by a finite collection of Balls \mathcal{F} . Show that then there exists a disjoint collection of these balls, $\{B_i\}_{i=1}^m$, such that $A \subseteq \cup_{i=1}^m \widehat{B}_i$ where \widehat{B}_i has the same center as B_i but 3 times the radius. **Hint:** Since the collection of balls is finite, they can be arranged in order of decreasing radius. Mimic the argument for Vitali covering theorem.

Chapter 5

Functions on Normed Linear Spaces

This chapter is about the general notion of functions defined on normed linear spaces even if the linear space is not finite dimensional.

5.1 $\mathcal{L}(V, W)$ as a Vector Space

In what follows, V, W will be vector spaces.

Definition 5.1.1 The term $\mathcal{L}(V, W)$ signifies the set of linear maps from V to W . This means that for $v, u \in V$ and α, β scalars from \mathbb{F} , $L(\alpha u + \beta v) = \alpha L(u) + \beta L(v)$. Given $L, M \in \mathcal{L}(V, W)$ define a new element of $\mathcal{L}(V, W)$, denoted by $L + M$ according to the rule¹ $(L + M)v \equiv Lv + Mv$. For α a scalar and $L \in \mathcal{L}(V, W)$, define $\alpha L \in \mathcal{L}(V, W)$ by $\alpha L(v) \equiv \alpha(Lv)$.

Note that if you have $V = \mathbb{R}^n$ and $W = \mathbb{R}^m$, an example of something in $\mathcal{L}(V, W)$ is given by $Tv \equiv Av$ where A is a real $m \times n$ matrix.

You should verify that all the axioms of a vector space hold for $\mathcal{L}(V, W)$ with the above definitions of vector addition and scalar multiplication. What about the dimension of $\mathcal{L}(V, W)$?

Before answering this question, here is a useful lemma. It gives a way to define linear transformations and a way to tell when two of them are equal.

Lemma 5.1.2 Let V and W be vector spaces and suppose $\{v_1, \dots, v_n\}$ is a basis for V . Then if $L: V \rightarrow W$ is given by $Lv_k = w_k \in W$ and $L(\sum_{k=1}^n a_k v_k) \equiv \sum_{k=1}^n a_k Lv_k = \sum_{k=1}^n a_k w_k$ then L is well defined and is in $\mathcal{L}(V, W)$. Also, if L, M are two linear transformations such that $Lv_k = Mv_k$ for all k , then $M = L$.

Proof: L is well defined on V because, since $\{v_1, \dots, v_n\}$ is a basis, there is exactly one way to write a given vector of V as a linear combination. Next, observe that L is obviously linear from the definition. If L, M are equal on the basis, then if $\sum_{k=1}^n a_k v_k$ is an arbitrary vector of V , $L(\sum_{k=1}^n a_k v_k) = \sum_{k=1}^n a_k Lv_k = \sum_{k=1}^n a_k Mv_k = M(\sum_{k=1}^n a_k v_k)$ and so $L = M$ because they give the same result for every vector in V . ■

The message is that when you define a linear transformation, it suffices to tell what it does to a basis.

Theorem 5.1.3 Let V and W be finite dimensional linear spaces of dimension n and m respectively. Then $\dim(\mathcal{L}(V, W)) = mn$.

Proof: Let two sets of bases be $\{v_1, \dots, v_n\}$ and $\{w_1, \dots, w_m\}$ for V and W respectively. Using Lemma 5.1.2, let $w_i v_j \in \mathcal{L}(V, W)$ be the linear transformation defined on the basis, $\{v_1, \dots, v_n\}$, by $w_i v_j(v_k) \equiv w_i \delta_{jk}$ where $\delta_{ik} = 1$ if $i = k$ and 0 if $i \neq k$. I will show that $L \in \mathcal{L}(V, W)$ is a linear combination of these special linear transformations called dyadics.

Then let $L \in \mathcal{L}(V, W)$. Since $\{w_1, \dots, w_m\}$ is a basis, there exist constants, d_{jk} such that $Lv_r = \sum_{j=1}^m d_{jr} w_j$. Now consider the following sum of dyadics. $\sum_{j=1}^m \sum_{i=1}^n d_{ji} w_j v_i$. Apply this to v_r . This yields $\sum_{j=1}^m \sum_{i=1}^n d_{ji} w_j v_i(v_r) = \sum_{j=1}^m \sum_{i=1}^n d_{ji} w_j \delta_{ir} = \sum_{j=1}^m d_{jr} w_j = Lv_r$. Therefore, $L = \sum_{j=1}^m \sum_{i=1}^n d_{ji} w_j v_i$ showing the span of the dyadics is all of $\mathcal{L}(V, W)$.

¹Note that this is the standard way of defining the sum of two functions.

Now consider whether these dyadics form a linearly independent set. Suppose that $\sum_{i,k} d_{ik} w_i v_k = \mathbf{0}$. Are all the scalars d_{ik} equal to 0? $\mathbf{0} = \sum_{i,k} d_{ik} w_i v_k(v_l) = \sum_{i=1}^m d_{il} w_i$ so, since $\{w_1, \dots, w_m\}$ is a basis, $d_{il} = 0$ for each $i = 1, \dots, m$. Since l is arbitrary, this shows $d_{il} = 0$ for all i and l . Thus these linear transformations form a basis and this shows that the dimension of $\mathcal{L}(V, W)$ is mn as claimed because there are m choices for the w_i and n choices for the v_j . ■

5.2 The Norm of a Linear Map, Operator Norm

Not surprisingly all of the above holds for a finite dimensional normed linear space. First here is an easy lemma which follows right away from Theorem 3.6.2, the theorem about equivalent formulations of continuity.

Lemma 5.2.1 *Let $(V, \|\cdot\|_V)$ and $(W, \|\cdot\|_W)$ be two normed linear spaces. Then a linear map $f : V \rightarrow W$ is continuous if and only if it takes bounded sets to bounded sets. (f is bounded) If V is finite dimensional, then f must be continuous.*

Proof: \Rightarrow Consider $f(B(0, 1))$. If this is not bounded, then there exists $\|v^m\|_V \leq 1$ but $\|f(v^m)\|_W \geq m$. Then it follows that $\left\|f\left(\frac{v^m}{m}\right)\right\|_W \geq 1$ which is impossible for all m since $\left\|\frac{v^m}{m}\right\|_V \leq \frac{1}{m}$ and so continuity requires that $\lim_{m \rightarrow \infty} f\left(\frac{v^m}{m}\right) = 0$ (Theorem 3.6.2). Thus there exists M such that $\|f(v)\| \leq M$ whenever $v \in B(0, 1)$. In general, let S be a bounded set. Then $S \subseteq B(0, r)$ for large enough r . Hence, for $v \in S$, it follows that $v/2r \in B(0, 1)$. It follows that $\|f(v/2r)\|_W \leq M$ and so $\|f(v)\|_W \leq 2rM$. Thus f takes bounded sets to bounded sets.

\Leftarrow Suppose f is bounded and not continuous. Then by Theorem 3.6.2 again, there is a sequence $v_n \rightarrow v$ but $f(v_n)$ fails to converge to $f(v)$. Then there exists $\varepsilon > 0$ and a subsequence, still denoted as v_n such that $\|f(v_n) - f(v)\| = \|f(v_n - v)\| \geq \varepsilon$. Then

$$\left\|f\left(\frac{v_n - v}{\|v_n - v\|}\right)\right\| \geq \varepsilon \frac{1}{\|v_n - v\|}$$

The right side is unbounded, but the left is bounded, a contradiction.

Consider the last claim about continuity. Let $\{v_1, \dots, v_n\}$ be a basis for V . By Lemma 4.4.7, if $y^m \rightarrow 0$, in V for $y^m = \sum_{k=1}^n y_k^m v_k$, then it follows that $\lim_{m \rightarrow \infty} y_k^m = 0$ and consequently, $f(y^m) \rightarrow f(0) = 0$. In general, if $y^m \rightarrow y$, then $(y^m - y) \rightarrow 0$ and so $f(y^m - y) = f(y^m) - f(y) \rightarrow 0$. That is, $f(y^m) \rightarrow f(y)$. ■

Definition 5.2.2 *For $f : (V, \|\cdot\|_V) \rightarrow (W, \|\cdot\|_W)$ continuous, it was just shown that there exists M such that $\|f(v)\| \leq M$, $v \in B(0, 1)$. It follows that, since $\frac{v}{\|v\|} \in B(0, 1)$, then $\|f(v)\| \leq 2M\|v\|$. Therefore, letting $\|f\| \equiv \sup_{\|v\| \leq 1} \|f(v)\|$ it follows that for all $v \in V$, $\|f(v)\| \leq \|f\| \|v\|$. Thus a linear map is bounded if and only if $\|f\| < \infty$ if and only if f is continuous. The number $\|f\|$ is called the operator norm. For X a real normed linear space, X' denotes the space $\mathcal{L}(X, \mathbb{R})$.*

You can show that for $\mathcal{L}(V, W)$ the space of bounded linear maps from V to W , $\mathcal{L}(V, W)$ becomes a normed linear space with this definition. This is true whether V, W are finite or infinite dimensional. You can also show that if W is complete then so is $\mathcal{L}(V, W)$. This is left as an exercise. Also, when the vector spaces are finite dimensional, Lemma 5.2.1 shows that any linear function f is automatically bounded, hence continuous, hence $\|f\|$ exists. Here is an interesting observation about the operator norm.

Lemma 5.2.3 Let $f \in \mathcal{L}(V, W)$ and let $h \in \mathcal{L}(W, Z)$ where X, Y, Z are normed vector spaces. Then $\|h \circ f\| \leq \|h\| \|f\|$.

Proof: This follows right away from the definition. If $\|v\| \leq 1$, then $\|f(v)\| \leq \|f\|$. This explains the first inequality in the following.

$$\sup_{\|v\| \leq 1} \|h \circ f(v)\| \leq \sup_{\|w\| \leq \|f\|} \|h(w)\| = \sup_{\|w\| \leq \|f\|} \left\| h \left(\frac{w}{\|f\|} \right) \right\| \|f\| \leq \|h\| \|f\|. \blacksquare$$

Theorem 5.2.4 Let $(V, \|\cdot\|)$ be a normed linear space with basis $\{v_1, \dots, v_n\}$ and field of scalars \mathbb{F} . Let $f: (\mathbb{F}^n, \|\cdot\|) \rightarrow (V, \|\cdot\|_V)$ be any linear map which is one to one and onto. Then both f and f^{-1} are continuous. Also the compact sets of $(V, \|\cdot\|_V)$ are exactly those which are closed and bounded.

Proof: Define another norm $\|\cdot\|_1$ on \mathbb{F}^n as follows. $\|x\|_1 \equiv \|f(x)\|_V$. Since f is one to one and onto and linear, this is indeed a norm. The details are left as an exercise. Then from the theorem on the equivalence of norms, there are positive constants δ, Δ such that $\delta \|x\| \leq \|f(x)\|_V \leq \Delta \|x\|$. Since f is one to one and onto, this implies $\delta \|f^{-1}(v)\| \leq \|v\|_V \leq \Delta \|f^{-1}(v)\|$. The first of these above inequalities implies f is continuous. The second says $\|f^{-1}(v)\| \leq \frac{1}{\delta} \|v\|_V$ and so f^{-1} is continuous. Thus, from the above theorems, both f and f^{-1} map closed sets to closed sets, compact sets to compact sets, open sets to open sets and bounded sets to bounded sets.

Now let $K \subseteq V$ be closed and bounded. Then from the above observations, $f^{-1}(K)$ is also closed and bounded. Therefore, it is compact. Now $f(f^{-1}(K)) = K$ must be compact because the continuous image of a compact set is compact, Theorem 3.7.1. Conversely, if $K \subseteq V$ is compact, then by the theorem just mentioned, $f^{-1}(K)$ is compact and so it is closed and bounded. Hence $f(f^{-1}(K)) = K$ is also closed and bounded. \blacksquare

This is a remarkable theorem. It says that an algebraic isomorphism is also a homeomorphism which is what it means to say that the map takes open sets to open sets and the inverse does the same. In other words, there really isn't any algebraic or topological distinction between a finite dimensional normed vector space of dimension n and \mathbb{F}^n . Of course when one considers geometry, this is not so.

Here is another interesting theorem about coordinate maps. It follows right away from earlier theorems.

Theorem 5.2.5 Let $f: (V, \|\cdot\|_V) \rightarrow (W, \|\cdot\|_W)$ be a continuous function where here $(V, \|\cdot\|_V)$ is a normed linear space and $(W, \|\cdot\|_W)$ is a finite dimensional normed linear space with basis $\{w_1, \dots, w_n\}$. Thus $f(v) \equiv \sum_{k=1}^n f_k(v) w_k$. Then f is continuous if and only if each f_k is a continuous \mathbb{F} valued map.

Proof: \implies First, why is f_k linear? This follows from

$$\begin{aligned} \sum_{k=1}^n (\alpha f_k(u) + \beta f_k(v)) w_k &= \alpha \sum_{k=1}^n f_k(u) w_k + \beta \sum_{k=1}^n f_k(v) w_k \\ &= \alpha f(u) + \beta f(v) = f(\alpha u + \beta v) \equiv \sum_{k=1}^n f_k(\alpha u + \beta v) w_k \end{aligned}$$

Why is the coordinate function f_k continuous? From Lemma 5.2.1, it suffices to verify that f_k is bounded. If this is not so, there exists $v_m, \|v_m\|_V \leq 1$ but $|f_k(v_m)|_W \geq m$. It follows

that $|f_k(\frac{v_m}{m})| \geq 1$. Since f is continuous, and $v_m/m \rightarrow 0$, it follows that $f(\frac{v_m}{m}) \rightarrow 0$ in V . However, by Lemma 4.4.7, $f_k(\frac{v_m}{m}) \rightarrow 0$, a contradiction.

\Leftarrow If each coordinate function is continuous, then

$$\|f(v) - f(\hat{v})\|_W = \left\| \sum_{k=1}^n f_k(v) w_k - \sum_{k=1}^n f_k(\hat{v}) w_k \right\| \leq \sum_{k=1}^n |f_k(v) - f_k(\hat{v})| \|w_k\|_W$$

Since each f_k is continuous, this shows that f is also. ■

5.3 Comparisons

Here are some useful lemmas about comparisons. Here $|\cdot|$ will be the usual norm but one could generalize.

Lemma 5.3.1 *Suppose S, T are linear, defined on a finite dimensional normed linear space, S^{-1} exists, and let $\delta \in (0, 1)$. Then whenever $\|S - T\|$ is small enough, it follows that*

$$\frac{|Tv|}{|Sv|} \in (1 - \delta, 1 + \delta) \quad (5.1)$$

for all $v \neq 0$. Similarly if T^{-1} exists and $\|S - T\|$ is small enough,

$$\frac{|Tv|}{|Sv|} \in (1 - \delta, 1 + \delta)$$

Proof: Say S^{-1} exists. Then $v \rightarrow |Sv|$ is a norm. Then by equivalence of norms, Theorem 4.4.9, there exists $\eta > 0$ such that for all v , $|Sv| \geq \eta |v|$. Say $\|T - S\| < r < \delta\eta$

$$\frac{|Tv|}{|Sv|} = \frac{|Sv - (S - T)v|}{|Sv|} \geq \frac{|Sv| - \|T - S\| |v|}{|Sv|} \geq \frac{|Sv| - \delta\eta |v|}{|Sv|} \geq \frac{|Sv| - \delta |Sv|}{|Sv|} = 1 - \delta$$

$$\frac{|Tv|}{|Sv|} = \frac{|Sv + (T - S)v|}{|Sv|} \leq \frac{|Sv| + \|T - S\| |v|}{|Sv|} \leq \frac{|Sv| + \delta\eta |v|}{|Sv|} \leq \frac{|Sv| + \delta |Sv|}{|Sv|} = 1 + \delta$$

Next suppose that T^{-1} exists. Then, letting $\hat{\delta}$ be small enough, $(1 - \hat{\delta}, 1 + \hat{\delta}) \subseteq (\frac{1}{1+\delta}, \frac{1}{1-\delta})$. From what was just shown, if $\|S - T\|$ is small enough,

$$\frac{|Sv|}{|Tv|} \in (1 - \hat{\delta}, 1 + \hat{\delta}) \subseteq \left(\frac{1}{1+\delta}, \frac{1}{1-\delta}\right) \text{ so } \frac{|Tv|}{|Sv|} \in (1 - \delta, 1 + \delta). \blacksquare$$

In short, the above lemma says that if one of S, T is invertible and the other is close to it, then it is also invertible and the quotient of $|Sv|$ and $|Tv|$ is close to 1. Then the following lemma is fairly obvious.

Lemma 5.3.2 *Let S, T be $n \times n$ matrices which are invertible. Then*

$$o(Tv) = o(Sv) = o(v)$$

and if L is a continuous linear transformation such that for $a < b$,

$$\sup_{v \neq 0} \frac{|Lv|}{|Sv|} < b, \quad \inf_{v \neq 0} \frac{|Lv|}{|Sv|} > a$$

If $\|S - T\|$ is small enough, it follows that the same inequalities hold with S replaced with T . Here $\|\cdot\|$ denotes the operator norm.

Proof: Consider the first claim. For

$$\frac{|o(Tv)|}{|v|} = \frac{|o(Tv)|}{|Tv|} \frac{|Tv|}{|v|} \leq \frac{|o(Tv)|}{|Tv|} \|T\|$$

Thus $o(Tv) = o(v)$. It is similar for T replaced with S .

Consider the second claim. Pick δ sufficiently small. Then by Lemma 17.2.1

$$\sup_{v \neq 0} \frac{|Lv|}{|Tv|} = \sup_{v \neq 0} \frac{|Lv|}{|Sv|} \frac{|Sv|}{|Tv|} \leq (1 + \delta) \sup_{v \neq 0} \frac{|Lv|}{|Sv|} < b$$

if δ is small enough. The other inequality is shown exactly similar. ■

5.4 Continuous Functions in Normed Linear Space

Of course not all functions are linear. Continuous functions have already been discussed in general metric space, but now there are other considerations to consider due to the algebra available in a normed linear space. The following theorem includes these kinds of considerations for functions having values in a normed linear space.

Theorem 5.4.1 *Let f, g be continuous functions defined on D , a metric space. Also let α, β be scalars. Then the following hold.*

1. $\alpha f + \beta g$ is continuous.
2. If $(W, \|\cdot\|_W)$ is an inner product space, then (f, g) defined as $(f, g)(v) \equiv (f(v), g(v))$, then (f, g) is continuous.
3. If f has values in \mathbb{F} and g has values in $(W, \|\cdot\|_W)$, then fg is continuous.

Proof: Say $v_n \rightarrow v$. Then

$$\|(\alpha f + \beta g)(v_n) - (\alpha f + \beta g)(v)\| \leq |\alpha| \|f(v_n) - f(v)\| + |\beta| \|g(v_n) - g(v)\|$$

and the right side converges to 0 as $n \rightarrow \infty$ so this shows 1.

This follows from an easy computation. From the Cauchy Schwarz inequality,

$$\begin{aligned} |(f, g)(v) - (f, g)(\hat{v})| &\leq |(f(v), g(v)) - (f(v), g(\hat{v}))| + |(f(v), g(\hat{v})) - (f(\hat{v}), g(\hat{v}))| \\ &\leq \|g(v) - g(\hat{v})\| \|f(v)\| + \|f(v) - f(\hat{v})\| \|g(\hat{v})\| \end{aligned}$$

Now since g is continuous at v and so $\|g(v) - g(\hat{v})\| < 1$ provided $d(v, \hat{v})$ is small enough. Thus $\|g(\hat{v})\| \leq \|g(v)\| + 1$. Hence if $d(v, \hat{v})$ is small enough,

$$|(f, g)(v) - (f, g)(\hat{v})| \leq (\|g(v)\| + 1) \|f(v) - f(\hat{v})\| + \|f(v)\| \|g(v) - g(\hat{v})\|$$

Thus, by continuity of f, g at v , if $d(v, \hat{v})$ is sufficiently small, the right side is less than ε and so $f \cdot g$ is continuous at v . This shows 2. The proof of 3. is just like this. ■

Of course there are other things like cross product and determinant and so forth which are defined in terms of the component functions of f . Then these things will be continuous by an application of Theorem 5.2.5.

5.5 Polynomials

For functions of one variable, the special kind of functions known as a polynomial has a corresponding version when one considers a function of many variables. This is found in the next definition.

Definition 5.5.1 *Let α be an n dimensional multi-index. The meaning of this term is that $\alpha = (\alpha_1, \dots, \alpha_n)$ where each α_i is a positive integer or zero. Also, let $|\alpha| \equiv \sum_{i=1}^n |\alpha_i|$. Then \mathbf{x}^α means $\mathbf{x}^\alpha \equiv x_1^{\alpha_1} x_2^{\alpha_2} \cdots x_n^{\alpha_n}$ where each $x_j \in \mathbb{F}$. An n dimensional polynomial of degree m is a function of the form $p(\mathbf{x}) = \sum_{|\alpha| \leq m} d_\alpha \mathbf{x}^\alpha$, where the d_α are complex or real numbers, more generally in some normed linear space X . Rational functions are defined as the quotient of two real or complex valued polynomials. Thus these functions are defined on \mathbb{F}^n .*

For example, $f(\mathbf{x}) = x_1 x_2^2 + 7x_3^4 x_1$ is a polynomial of degree 5 and $\frac{x_1 x_2^2 + 7x_3^4 x_1 + x_2^3}{4x_1^3 x_2^2 + 7x_3^2 x_1 - x_2^3}$ is a rational function.

Note that in the case of a rational function, the domain of the function might not be all of \mathbb{F}^n . For example, if $f(\mathbf{x}) = \frac{x_1 x_2^2 + 7x_3^4 x_1 + x_2^3}{x_2^2 + 3x_1^2 - 4}$, the domain of f would be all complex numbers such that $x_2^2 + 3x_1^2 \neq 4$.

By Theorem 3.6.2 all polynomials are continuous. To see this, note that the function, $\pi_k(\mathbf{x}) \equiv x_k$ is a continuous function because of the inequality

$$|\pi_k(\mathbf{x}) - \pi_k(\mathbf{y})| = |x_k - y_k| \leq |\mathbf{x} - \mathbf{y}|.$$

Polynomials are simple sums of scalars times products of these functions. Similarly, by this theorem, rational functions, quotients of polynomials, are continuous at points where the denominator is non zero. More generally, if V is a normed vector space, consider a V valued function of the form $\mathbf{f}(\mathbf{x}) \equiv \sum_{|\alpha| \leq m} \mathbf{d}_\alpha \mathbf{x}^\alpha$ where $\mathbf{d}_\alpha \in V$, sort of a V valued polynomial. Then such a function is continuous by application of Theorem 3.6.2 and the above observation about the continuity of the functions π_k .

Thus there are lots of examples of continuous functions. However, it is even better than the above discussion indicates. As in the case of a function of one variable, an arbitrary continuous function can typically be approximated uniformly by a polynomial. This is the n dimensional version of the Weierstrass approximation theorem.

5.6 Weierstrass Approximation Theorem

An arbitrary continuous function defined on an interval can be approximated uniformly by a polynomial, there exists a similar theorem which is just a generalization of this which will hold for continuous functions defined on a box or more generally a closed and bounded set. However, we will settle for the case of a box first. The proof is based on the following lemma.

Lemma 5.6.1 *The following estimate holds for $x \in [0, 1]$ and $m \geq 2$.*

$$\sum_{k=0}^m \binom{m}{k} (k - mx)^2 x^k (1 - x)^{m-k} \leq \frac{1}{4} m$$

Proof: First of all, from the binomial theorem,

$$\begin{aligned} \sum_{k=0}^m \binom{m}{k} \left(e^{t(k-mx)} \right) x^k (1-x)^{m-k} &= e^{-tmx} \sum_{k=0}^m \binom{m}{k} \left(e^{tk} \right) x^k (1-x)^{m-k} \\ &= e^{-tmx} (1-x+xe^t)^m = e^{-tmx} g(t)^m, \quad g(0) = 1, g'(0) = g''(0) = x \end{aligned}$$

Take a partial derivative with respect to t twice.

$$\begin{aligned} &\sum_{k=0}^m \binom{m}{k} (k-mx)^2 e^{t(k-mx)} x^k (1-x)^{m-k} \\ &= (mx)^2 e^{-tmx} g(t)^m + 2(-mx) e^{-tmx} m g(t)^{m-1} g'(t) \\ &\quad + e^{-tmx} \left[m(m-1) g(t)^{m-2} g'(t)^2 + m g(t)^{m-1} g''(t) \right] \end{aligned}$$

Now let $t = 0$ and note that the right side is $m(x-x^2) \leq m/4$ for $x \in [0, 1]$. Thus

$$\sum_{k=0}^m \binom{m}{k} (k-mx)^2 x^k (1-x)^{m-k} = mx - mx^2 \leq m/4 \blacksquare$$

With this preparation, here is the first version of the Weierstrass approximation theorem. I will allow f to have values in a complete, real or complex normed linear space. Thus, $f \in C([0, 1]; X)$ where X is a Banach space, Definition 4.3.7. Thus this is a function which is continuous with values in X as discussed earlier with metric spaces.

Theorem 5.6.2 *Let $f \in C([0, 1]; X)$ and let the norm on X be denoted by $\|\cdot\|$.*

$$p_m(x) \equiv \sum_{k=0}^m \binom{m}{k} x^k (1-x)^{m-k} f\left(\frac{k}{m}\right) = \sum_{k=0}^m q_k(x) f\left(\frac{k}{m}\right)$$

Then these polynomials having coefficients in X converge uniformly to f on $[0, 1]$. Also $q_0(0) = 1, q_k(0) = 0$ for $k \neq 0$, and $q_m(1) = 1$ while $q_k(1) = 0$ for $k \neq m$.

Proof: Let $\|f\|_\infty$ denote the largest value of $\|f(x)\|$. By uniform continuity of f , there exists a $\delta > 0$ such that if $|x - x'| < \delta$, then $\|f(x) - f(x')\| < \varepsilon/2$. By the binomial theorem,

$$\begin{aligned} \|p_m(x) - f(x)\| &\leq \sum_{k=0}^m \binom{m}{k} x^k (1-x)^{m-k} \left\| f\left(\frac{k}{m}\right) - f(x) \right\| \\ &\leq \sum_{\left|\frac{k}{m} - x\right| < \delta} \binom{m}{k} x^k (1-x)^{m-k} \left\| f\left(\frac{k}{m}\right) - f(x) \right\| + 2\|f\|_\infty \sum_{\left|\frac{k}{m} - x\right| \geq \delta} \binom{m}{k} x^k (1-x)^{m-k} \end{aligned}$$

Therefore,

$$\begin{aligned} &\leq \sum_{k=0}^m \binom{m}{k} x^k (1-x)^{m-k} \frac{\varepsilon}{2} + 2\|f\|_\infty \sum_{(k-mx)^2 \geq m^2 \delta^2} \binom{m}{k} x^k (1-x)^{m-k} \\ &\leq \frac{\varepsilon}{2} + 2\|f\|_\infty \frac{1}{m^2 \delta^2} \sum_{k=0}^m \binom{m}{k} (k-mx)^2 x^k (1-x)^{m-k} \leq \frac{\varepsilon}{2} + 2\|f\|_\infty \frac{1}{4} m \frac{1}{\delta^2 m^2} < \varepsilon \end{aligned}$$

provided m is large enough. Thus $\|p_m - f\|_\infty < \varepsilon$ when m is large enough. \blacksquare

Note that we do not need to have X be complete in order for this to hold. It would have sufficed to have simply let X be a normed linear space.

Corollary 5.6.3 *If $f \in C([a, b]; X)$ where X is a normed linear space, then there exists a sequence of polynomials which converge uniformly to f on $[a, b]$. The m^{th} term of this sequence is $\sum_{k=0}^m q_k(y) f\left(l\left(\frac{k}{m}\right)\right)$ where $l: [0, 1] \rightarrow [a, b]$ be one to one, linear and onto and $q_0(a) = 1$ and if $k \neq 0, q_k(a) = 0$ and $q_m(b) = 1$ and if $k \neq m$, then $q_k(b) = 0$.*

Proof: Let $l: [0, 1] \rightarrow [a, b]$ be one to one, linear and onto. Then $f \circ l$ is continuous on $[0, 1]$ and so if $\varepsilon > 0$ is given, if m large enough, then for all $x \in [0, 1]$,

$$\left\| \sum_{k=0}^m \hat{q}_k(x) f\left(l\left(\frac{k}{m}\right)\right) - f \circ l(x) \right\| < \varepsilon$$

where $\hat{q}_0(0) = 1$ and $\hat{q}_k(0) = 0$ for $k \neq 0, \hat{q}_m(1) = 1, \hat{q}_k(1) = 0$ if $k \neq m$. Therefore, for all $y \in [a, b]$,

$$\left\| \sum_{k=0}^m \hat{q}_k(l^{-1}(y)) f\left(l\left(\frac{k}{m}\right)\right) - f(y) \right\| < \varepsilon$$

Let $q_k(y) \equiv \hat{q}_k(l^{-1}(y))$. ■

As another corollary, here is the version which will be used in Stone's generalization later.

Corollary 5.6.4 *Let f be a continuous function defined on $[-M, M]$ with $f(0) = 0$. Then there is a sequence of polynomials $\{p_m\}$, $p_m(0) = 0$ and*

$$\lim_{m \rightarrow \infty} \|p_m - f\|_{\infty} = 0$$

Proof: From Corollary 5.6.3 there exists a sequence of polynomials $\{\widehat{p}_m\}$ such that $\|\widehat{p}_m - f\|_{\infty} \rightarrow 0$. Simply consider $p_m = \widehat{p}_m - \widehat{p}_m(0)$. ■

5.7 Functions of Many Variables

First note that if $h: K \times H \rightarrow \mathbb{R}$ is a real valued continuous function where K, H are compact sets in metric spaces,

$$\max_{x \in K} h(x, y) \geq h(x, y), \text{ so } \max_{y \in H} \max_{x \in K} h(x, y) \geq h(x, y)$$

which implies $\max_{y \in H} \max_{x \in K} h(x, y) \geq \max_{(x, y) \in K \times H} h(x, y)$. The other inequality is also obtained.

Let $f \in C(R_p; X)$ where $R_p = [0, 1]^p$. Then let $\hat{x}_p \equiv (x_1, \dots, x_{p-1})$. By Theorem 5.6.2, if n is large enough,

$$\max_{x_p \in [0, 1]} \left\| \sum_{k=0}^n f\left(\cdot, \frac{k}{n}\right) \binom{n}{k} x_p^k (1 - x_p)^{n-k} - f(\cdot, x_p) \right\|_{C([0, 1]^{p-1}; X)} < \frac{\varepsilon}{2}$$

Now $f(\cdot, \frac{k}{n}) \in C(R_{p-1}; X)$ and so by induction, there is a polynomial $p_k(\hat{x}_p)$ such that

$$\max_{\hat{x}_p \in R_{p-1}} \left\| p_k(\hat{x}_p) - \binom{n}{k} f\left(\hat{x}_p, \frac{k}{n}\right) \right\|_X < \frac{\varepsilon}{(n+1)2}$$

Thus, letting $\mathbf{p}(\mathbf{x}) \equiv \sum_{k=0}^n \mathbf{p}_k(\hat{\mathbf{x}}_p) x_p^k (1-x_p)^{n-k}$,

$$\|\mathbf{p} - \mathbf{f}\|_{C(R_p; X)} \leq \max_{x_p \in [0,1]} \max_{\hat{\mathbf{x}}_p \in R_{p-1}} \|\mathbf{p}(\hat{\mathbf{x}}_p, x_p) - \mathbf{f}(\hat{\mathbf{x}}_p, x_p)\|_X < \varepsilon$$

where \mathbf{p} is a polynomial with coefficients in X .

In general, if $R_p \equiv \prod_{k=1}^p [a_k, b_k]$, note that there is a linear function $l_k : [0, 1] \rightarrow [a_k, b_k]$ which is one to one and onto. Thus $\mathbf{l}(\mathbf{x}) \equiv (l_1(x_1), \dots, l_p(x_p))$ is a one to one and onto map from $[0, 1]^p$ to R_p and the above result can be applied to $\mathbf{f} \circ \mathbf{l}$ to obtain a polynomial \mathbf{p} with $\|\mathbf{p} - \mathbf{f} \circ \mathbf{l}\|_{C([0,1]^p; X)} < \varepsilon$. Thus $\|\mathbf{p} \circ \mathbf{l}^{-1} - \mathbf{f}\|_{C(R_p; X)} < \varepsilon$ and $\mathbf{p} \circ \mathbf{l}^{-1}$ is a polynomial. This proves the following theorem.

Theorem 5.7.1 *Let \mathbf{f} be a function in $C(R; X)$ for X a normed linear space where $R \equiv \prod_{k=1}^p [a_k, b_k]$. Then for any $\varepsilon > 0$ there exists a polynomial \mathbf{p} having coefficients in X such that $\|\mathbf{p} - \mathbf{f}\|_{C(R; X)} < \varepsilon$.*

These Bernstein polynomials are very remarkable approximations. It turns out that if f is $C^1([0, 1]; X)$, then $\lim_{n \rightarrow \infty} p'_n(x) \rightarrow f'(x)$ uniformly on $[0, 1]$. This all works for functions of many variables as well, but here I will only show it for functions of one variable.

Lemma 5.7.2 *Let $f \in C^1([0, 1])$ and let $p_m(x) \equiv \sum_{k=0}^m \binom{m}{k} x^k (1-x)^{m-k} f\left(\frac{k}{m}\right)$ be the m^{th} Bernstein polynomial. Then in addition to $\|p_m - f\|_{[0,1]} \rightarrow 0$, it also follows that $\|p'_m - f'\|_{[0,1]} \rightarrow 0$.*

Proof: From simple computations,

$$\begin{aligned} p'_m(x) &= \sum_{k=1}^m \binom{m}{k} k x^{k-1} (1-x)^{m-k} f\left(\frac{k}{m}\right) \\ &\quad - \sum_{k=0}^{m-1} \binom{m}{k} x^k (m-k) (1-x)^{m-1-k} f\left(\frac{k}{m}\right) \\ &= \sum_{k=1}^m \frac{m(m-1)!}{(m-k)!(k-1)!} x^{k-1} (1-x)^{m-k} f\left(\frac{k}{m}\right) \\ &\quad - \sum_{k=0}^{m-1} \binom{m}{k} x^k (m-k) (1-x)^{m-1-k} f\left(\frac{k}{m}\right) \\ &= \sum_{k=0}^{m-1} \frac{m(m-1)!}{(m-1-k)!k!} x^k (1-x)^{m-1-k} f\left(\frac{k+1}{m}\right) \\ &\quad - \sum_{k=0}^{m-1} \frac{m(m-1)!}{(m-1-k)!k!} x^k (1-x)^{m-1-k} f\left(\frac{k}{m}\right) \\ &= \sum_{k=0}^{m-1} \frac{m(m-1)!}{(m-1-k)!k!} x^k (1-x)^{m-1-k} \left(f\left(\frac{k+1}{m}\right) - f\left(\frac{k}{m}\right) \right) \end{aligned}$$

$$= \sum_{k=0}^{m-1} \binom{m-1}{k} x^k (1-x)^{m-1-k} \left(\frac{f\left(\frac{k+1}{m}\right) - f\left(\frac{k}{m}\right)}{1/m} \right)$$

By the mean value theorem, $\frac{f\left(\frac{k+1}{m}\right) - f\left(\frac{k}{m}\right)}{1/m} = f'(x_{k,m})$, $x_{k,m} \in \left(\frac{k}{m}, \frac{k+1}{m}\right)$. Now the desired result follows as before from the uniform continuity of f' on $[0, 1]$. Let $\delta > 0$ be such that if $|x - y| < \delta$, then $|f'(x) - f'(y)| < \varepsilon$ and let m be so large that $1/m < \delta/2$. Then if $|x - \frac{k}{m}| < \delta/2$, it follows that $|x - x_{k,m}| < \delta$ and so

$$|f'(x) - f'(x_{k,m})| = \left| f'(x) - \frac{f\left(\frac{k+1}{m}\right) - f\left(\frac{k}{m}\right)}{1/m} \right| < \varepsilon.$$

Now as before, letting $M \geq |f'(x)|$ for all x ,

$$\begin{aligned} |p'_m(x) - f'(x)| &\leq \sum_{k=0}^{m-1} \binom{m-1}{k} x^k (1-x)^{m-1-k} |f'(x_{k,m}) - f'(x)| \\ &\leq \sum_{\{x: |x - \frac{k}{m}| < \frac{\delta}{2}\}} \binom{m-1}{k} x^k (1-x)^{m-1-k} \varepsilon \\ &\quad + M \sum_{k=0}^{m-1} \binom{m-1}{k} \frac{4(k-mx)^2}{m^2 \delta^2} x^k (1-x)^{m-1-k} \\ &\leq \varepsilon + 4M \frac{1}{4} m \frac{1}{m^2 \delta^2} = \varepsilon + M \frac{1}{m \delta^2} < 2\varepsilon \end{aligned}$$

whenever m is large enough. Thus this proves uniform convergence. ■

There is a more general version of the Weierstrass theorem which is easy to get. It depends on the Tietze extension theorem, a wonderful little result which is interesting for its own sake.

5.8 A Generalization

This is an interesting theorem which holds in arbitrary normal topological spaces. In particular it holds in metric space and this is the context in which it will be discussed. First, review Lemma 3.12.1.

Lemma 5.8.1 *Let H, K be two nonempty disjoint closed subsets of X . Then there exists a continuous function, $g : X \rightarrow [-1/3, 1/3]$ such that $g(H) = -1/3$, $g(K) = 1/3$, $g(X) \subseteq [-1/3, 1/3]$.*

Proof: Let $f(x) \equiv \frac{\text{dist}(x, H)}{\text{dist}(x, H) + \text{dist}(x, K)}$. The denominator is never equal to zero because if $\text{dist}(x, H) = 0$, then $x \in H$ because H is closed. (To see this, pick $h_k \in B(x, 1/k) \cap H$. Then $h_k \rightarrow x$ and since H is closed, $x \in H$.) Similarly, if $\text{dist}(x, K) = 0$, then $x \in K$ and so the denominator is never zero as claimed. Hence f is continuous and from its definition, $f = 0$ on H and $f = 1$ on K . Now let $g(x) \equiv \frac{2}{3} \left(f(x) - \frac{1}{2} \right)$. Then g has the desired properties. ■

Definition 5.8.2 For $f : M \subseteq X \rightarrow \mathbb{R}$, let $\|f\|_M \equiv \sup \{|f(x)| : x \in M\}$. This is just notation. I am not claiming this is a norm.

Lemma 5.8.3 Suppose M is a closed set in X and suppose $f : M \rightarrow [-1, 1]$ is continuous at every point of M . Then there exists a function, g which is defined and continuous on all of X such that $\|f - g\|_M \leq \frac{2}{3}$, $g(X) \subseteq [-1/3, 1/3]$. If X is a normed vector space, and f is odd, meaning that M is symmetric ($x \in M$ if and only if $-x \in M$) and $f(-x) = -f(x)$. Then we can assume g is also odd.

Proof: Let $H = f^{-1}([-1, -1/3])$, $K = f^{-1}([1/3, 1])$. Thus H and K are disjoint closed subsets of M . Suppose first H, K are both nonempty. Then by Lemma 5.8.1 there exists g such that g is a continuous function defined on all of X and $g(H) = -1/3$, $g(K) = 1/3$, and $g(X) \subseteq [-1/3, 1/3]$. It follows $\|f - g\|_M < 2/3$. If $H = \emptyset$, then f has all its values in $[-1/3, 1]$ and so letting $g \equiv 1/3$, the desired condition is obtained. If $K = \emptyset$, let $g \equiv -1/3$. If both $H, K = \emptyset$, let $g = 0$.

When M is symmetric and f is odd, $g(x) \equiv \frac{1}{3} \frac{\text{dist}(x, H) - \text{dist}(x, K)}{\text{dist}(x, H) + \text{dist}(x, K)}$. When $x \in H$ this gives $\frac{1 - \text{dist}(x, K)}{3 \text{dist}(x, K)} = -\frac{1}{3}$. Then $x \in K$, this gives $\frac{1}{3} \frac{\text{dist}(x, H)}{\text{dist}(x, H)} = \frac{1}{3}$. Also $g(H) = -1/3$, $f(H) \subseteq [-1, -1/3]$ so for $x \in H$, $|g(x) - f(x)| \leq \frac{2}{3}$. It is similar for $x \in K$. If x is in neither H nor K , then $g(x) \in [-1/3, 1/3]$ and so is $f(x)$. Thus $\|f - g\|_M \leq \frac{2}{3}$. Now by assumption, since f is odd, $H = -K$. It is clear that g is odd because

$$\begin{aligned} g(-x) &= \frac{1}{3} \frac{\text{dist}(-x, H) - \text{dist}(-x, K)}{\text{dist}(-x, H) + \text{dist}(-x, K)} = \frac{1}{3} \frac{\text{dist}(-x, -K) - \text{dist}(-x, -H)}{\text{dist}(-x, -K) + \text{dist}(-x, -H)} \\ &= \frac{1}{3} \frac{\text{dist}(x, K) - \text{dist}(x, H)}{\text{dist}(x, K) + \text{dist}(x, H)} = -g(x). \quad \blacksquare \end{aligned}$$

Lemma 5.8.4 Suppose M is a closed set in X and suppose $f : M \rightarrow [-1, 1]$ is continuous at every point of M . Then there exists a function g which is defined and continuous on all of X such that $g = f$ on M and g has its values in $[-1, 1]$. If X is a normed linear space and f is odd, then we can also assume g is odd.

Proof: Using Lemma 5.8.3, let g_1 be such that $g_1(X) \subseteq [-1/3, 1/3]$ and $\|f - g_1\|_M \leq \frac{2}{3}$. Suppose g_1, \dots, g_m have been chosen such that $g_j(X) \subseteq [-1/3, 1/3]$ and

$$\left\| f - \sum_{i=1}^m \left(\frac{2}{3}\right)^{i-1} g_i \right\|_M < \left(\frac{2}{3}\right)^m. \quad (5.2)$$

This has been done for $m = 1$. Then $\left\| \left(\frac{3}{2}\right)^m \left(f - \sum_{i=1}^m \left(\frac{2}{3}\right)^{i-1} g_i \right) \right\|_M \leq 1$ and so

$$\left(\frac{3}{2}\right)^m \left(f - \sum_{i=1}^m \left(\frac{2}{3}\right)^{i-1} g_i \right)$$

can play the role of f in the first step of the proof. Therefore, there exists g_{m+1} defined and continuous on all of X such that its values are in $[-1/3, 1/3]$ and

$$\left\| \left(\frac{3}{2}\right)^m \left(f - \sum_{i=1}^m \left(\frac{2}{3}\right)^{i-1} g_i \right) - g_{m+1} \right\|_M \leq \frac{2}{3}.$$

Hence

$$\left\| \left(f - \sum_{i=1}^m \left(\frac{2}{3} \right)^{i-1} g_i \right) - \left(\frac{2}{3} \right)^m g_{m+1} \right\|_M \leq \left(\frac{2}{3} \right)^{m+1}.$$

It follows there exists a sequence, $\{g_i\}$ such that each has its values in $[-1/3, 1/3]$ and for every m 5.2 holds. Then let $g(x) \equiv \sum_{i=1}^{\infty} \left(\frac{2}{3} \right)^{i-1} g_i(x)$. It follows

$$|g(x)| \leq \left| \sum_{i=1}^{\infty} \left(\frac{2}{3} \right)^{i-1} g_i(x) \right| \leq \sum_{i=1}^m \left(\frac{2}{3} \right)^{i-1} \frac{1}{3} \leq 1$$

and $\left| \left(\frac{2}{3} \right)^{i-1} g_i(x) \right| \leq \left(\frac{2}{3} \right)^{i-1} \frac{1}{3}$ so the Weierstrass M test applies and shows convergence is uniform. Therefore g must be continuous by Theorem 3.9.3. The estimate 5.2 implies $f = g$ on M . The last claim follows because we can take each g_i odd. ■

The following is the Tietze extension theorem.

Theorem 5.8.5 *Let M be a closed nonempty subset of a metric space X and let $f : M \rightarrow [a, b]$ be continuous at every point of M . Then there exists a function g continuous on all of X which coincides with f on M such that $g(X) \subseteq [a, b]$. If $[a, b]$ is centered on 0, and if X is a normed linear space and f is odd, then we can obtain that g is also odd.*

Proof: Let $f_1(x) = 1 + \frac{2}{b-a}(f(x) - b)$. Then f_1 satisfies the conditions of Lemma 5.8.4 and so there exists $g_1 : X \rightarrow [-1, 1]$ such that g is continuous on X and equals f_1 on M . Let $g(x) = (g_1(x) - 1) \left(\frac{b-a}{2} \right) + b$. This works. The last claim follows from the same arguments which gave Lemma 5.8.4 or the change of variables just given. ■

Corollary 5.8.6 *Let M be a closed nonempty subset of a metric space X and let $f : M \rightarrow [a, b]$ be continuous at every point of M . Also let $\|f - g\| \leq \varepsilon$. Then there exists continuous \hat{f} extending f with $\hat{f}(X) \subseteq [a, b]$ and \hat{g} extending g such that $\hat{g}(X) \subseteq [a - \varepsilon, b + \varepsilon]$. Also $\|\hat{f} - \hat{g}\| \leq \varepsilon$.*

Proof: Let \hat{f} be the extension of f from the above theorem. Now let F be the extension of $f - g$ with $\|F\| \leq \varepsilon$. Then let $\hat{g} = \hat{f} - F$. Then for $x \in M$, $\hat{g}(x) = f(x) - (f(x) - g(x)) = g(x)$. Thus it extends g and clearly $\hat{g}(X) \subseteq [a - \varepsilon, b + \varepsilon]$. ■

With the Tietze extension theorem, here is a better version of the Weierstrass approximation theorem.

Theorem 5.8.7 *Let K be a closed and bounded subset of \mathbb{R}^p and let $f : K \rightarrow \mathbb{R}$ be continuous. Then there exists a sequence of polynomials $\{p_m\}$ such that*

$$\lim_{m \rightarrow \infty} (\sup \{|f(x) - p_m(x)| : x \in K\}) = 0.$$

In other words, the sequence of polynomials converges uniformly to f on K .

Proof: By the Tietze extension theorem, there exists an extension of f to a continuous function g defined on all \mathbb{R}^p such that $g = f$ on K . Now since K is bounded, there exist intervals, $[a_k, b_k]$ such that $K \subseteq \prod_{k=1}^p [a_k, b_k] = R$. Then by the Weierstrass approximation theorem, Theorem 5.7.1 there exists a sequence of polynomials $\{p_m\}$ converging uniformly to g on R . Therefore, this sequence of polynomials converges uniformly to $g = f$ on K as well. This proves the theorem. ■

By considering the real and imaginary parts of a function which has values in \mathbb{C} one can generalize the above theorem.

Corollary 5.8.8 *Let K be a closed and bounded subset of \mathbb{R}^p and let $f : K \rightarrow \mathbb{F}$ be continuous. Then there exists a sequence of polynomials $\{p_m\}$ such that*

$$\lim_{m \rightarrow \infty} (\sup \{|f(x) - p_m(x)| : x \in K\}) = 0.$$

In other words, the sequence of polynomials converges uniformly to f on K .

More generally, the function f could have values in \mathbb{R}^p . There is no change in the proof. You just use norm symbols rather than absolute values and nothing at all changes in the theorem where the function is defined on a rectangle. Then you apply the Tietze extension theorem to each component in the case the function has values in \mathbb{R}^p .

5.9 An Approach to the Integral

With the Weierstrass approximation theorem, you can give a rigorous definition of the Riemann integral without wading in to Riemann sums. This shows the integral can be defined directly from very simple ideas. First is a short review of the derivative of a function of one variable.

Definition 5.9.1 *Let $f : [a, b] \rightarrow \mathbb{R}$. Then $f'(x) \equiv \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$ where h is always such that $x, x+h$ are both in the interval $[a, b]$ so we include derivatives at the right and left end points in this definition.*

The most important theorem about derivatives of functions of one variable is the mean value theorem.

Theorem 5.9.2 *Let $f : [a, b] \rightarrow \mathbb{R}$ be continuous. Then if the maximum value of f occurs at a point $x \in (a, b)$, it follows that if $f'(x) \neq 0$. If f achieves a minimum at $x \in (a, b)$ where $f'(x)$ exists, it also follows that $f'(x) = 0$.*

Proof: By Theorem 3.7.2, f achieves a maximum at some point x . If $f'(x)$ exists, then

$$f'(x) = \lim_{h \rightarrow 0+} \frac{f(x+h) - f(x)}{h} = \lim_{h \rightarrow 0-} \frac{f(x+h) - f(x)}{h}$$

However, the first limit is non-positive while the second is non-negative and so $f'(x) = 0$. The situation is similar if the minimum occurs at $x \in (a, b)$. ■

The Cauchy mean value theorem follows. The usual one is obtained by letting $g(x) = x$.

Theorem 5.9.3 *Let f, g be continuous on $[a, b]$ and differentiable on (a, b) . Then there exists $x \in (a, b)$ such that $f'(x)(g(b) - g(a)) = g'(x)(f(b) - f(a))$. If $g(x) = x$, this yields $f(b) - f(a) = f'(x)(b - a)$, also $f(a) - f(b) = f'(x)(a - b)$.*

Proof: Let $h(x) \equiv f(x)(g(b) - g(a)) - g(x)(f(b) - f(a))$. Then

$$h(a) = h(b) = f(a)g(b) - g(a)f(b).$$

If h is constant, then pick any $x \in (a, b)$ and $h'(x) = 0$. If h is not constant, then it has either a maximum or a minimum on (a, b) and so if x is the point where either occurs, then $h'(x) = 0$ which proves the theorem. ■

Recall that an antiderivative of a function f is just a function F such that $F' = f$. You know how to find an antiderivative for a polynomial. $\left(\frac{x^{n+1}}{n+1}\right)' = x^n$ so $\int \sum_{k=1}^n a_k x^k = \sum_{k=1}^n a_k \frac{x^{k+1}}{k+1} + C$. With this information and the Weierstrass theorem, it is easy to define integrals of continuous functions with all the properties presented in elementary calculus courses. It is an approach which does not depend on Riemann sums yet still gives the fundamental theorem of calculus. Note that if $F'(x) = 0$ for x in an interval, then for x, y in that interval, $F(y) - F(x) = 0(y - x)$ so F is a constant. Thus, if $F' = G'$ on an open interval, F, G continuous on the closed interval, it follows that $F - G$ is a constant and so $F(b) - F(a) = G(b) - G(a)$.

Definition 5.9.4 For $p(x)$ a polynomial on $[a, b]$, let $P'(x) = p(x)$. Thus, by the mean value theorem if P', \hat{P}' both equal p , it follows that $P(b) - P(a) = \hat{P}(b) - \hat{P}(a)$. Then define $\int_a^b p(x) dx \equiv P(b) - P(a)$. If $f \in C([a, b])$, define $\int_a^b f(x) dx \equiv \lim_{n \rightarrow \infty} \int_a^b p_n(x) dx$ where

$$\lim_{n \rightarrow \infty} \|p_n - f\| \equiv \lim_{n \rightarrow \infty} \max_{x \in [a, b]} |f(x) - p_n(x)| = 0$$

Proposition 5.9.5 The above integral is well defined and satisfies the following properties.

1. $\int_a^b f dx = f(\hat{x})(b - a)$ for some \hat{x} between a and b . Thus $\left| \int_a^b f dx \right| \leq \|f\| |b - a|$.
2. If f is continuous on an interval which contains all necessary intervals,

$$\int_a^c f dx + \int_c^b f dx = \int_a^b f dx, \text{ so } \int_a^b f dx + \int_b^a f dx = \int_b^b f dx = 0$$

3. If $F(x) \equiv \int_a^x f dt$, Then $F'(x) = f(x)$ so any continuous function has an antiderivative, and for any $a \neq b$, $\int_a^b f dx = G(b) - G(a)$ whenever $G' = f$ on the open interval determined by a, b and G continuous on the closed interval determined by a, b . Also,

$$\int_a^b (\alpha f(x) + \beta g(x)) dx = \alpha \int_a^b f(x) dx + \beta \int_a^b g(x) dx$$

If $a < b$, and $f(x) \geq 0$, then $\int_a^b f dx \geq 0$. Also $\left| \int_a^b f dx \right| \leq \left| \int_a^b |f| dx \right|$.

4. $\int_a^b 1 dx = b - a$.

Proof: First, why is the integral well defined? With notation as in the above definition, the mean value theorem implies

$$\int_a^b p(x) dx \equiv P(b) - P(a) = p(\hat{x})(b - a) \quad (5.3)$$

where \hat{x} is between a and b and so $\left| \int_a^b p(x) dx \right| \leq \|p\| |b - a|$. If $\|p_n - f\| \rightarrow 0$, then $\lim_{m, n \rightarrow \infty} \|p_n - p_m\| = 0$ and so

$$\left| \int_a^b p_n(x) dx - \int_a^b p_m(x) dx \right| = |(P_n(b) - P_n(a)) - (P_m(b) - P_m(a))|$$

$$= |(P_n(b) - P_m(b)) - (P_n(a) - P_m(a))| = \left| \int_a^b (p_n - p_m) dx \right| \leq \|p_n - p_m\| |b - a|$$

Thus the limit exists because $\left\{ \int_a^b p_n dx \right\}_n$ is a Cauchy sequence and \mathbb{R} is complete.

From 5.3, 1. holds for a polynomial $p(x)$. Let $\|p_n - f\| \rightarrow 0$. Then by definition,

$$\int_a^b f dx \equiv \lim_{n \rightarrow \infty} \int_a^b p_n dx = p_n(x_n)(b - a) \quad (5.4)$$

for some x_n in the open interval determined by (a, b) . By compactness, there is a further subsequence, still denoted with n such that $x_n \rightarrow x \in [a, b]$. Then fixing m such that $\|f - p_m\| < \varepsilon$ whenever $n \geq m$, assume $n > m$. Then $\|p_m - p_n\| \leq \|p_m - f\| + \|f - p_n\| < 2\varepsilon$ and so

$$\begin{aligned} |f(x) - p_n(x_n)| &\leq |f(x) - f(x_n)| + |f(x_n) - p_m(x_n)| + |p_m(x_n) - p_n(x_n)| \\ &\leq |f(x) - f(x_n)| + \|f - p_m\| + \|p_m - p_n\| < |f(x) - f(x_n)| + 3\varepsilon \end{aligned}$$

Now if n is still larger, continuity of f shows that $|f(x) - p_n(x_n)| < 4\varepsilon$. Since ε is arbitrary, $p_n(x_n) \rightarrow f(x)$ and so, passing to the limit with this subsequence in 5.4 yields 1.

Now consider 2. It holds for polynomials $p(x)$ obviously. So let $\|p_n - f\| \rightarrow 0$. Then

$$\int_a^c p_n dx + \int_c^b p_n dx = \int_a^b p_n dx$$

Pass to a limit as $n \rightarrow \infty$ and use the definition to get 2. Also note that $\int_b^b f(x) dx = 0$ follows from the definition.

Next consider 3. Let $h \neq 0$ and let x be in the open interval determined by a and b . Then for small h ,

$$\frac{F(x+h) - F(x)}{h} = \frac{1}{h} \int_x^{x+h} f(t) dt = f(x_h)$$

where x_h is between x and $x+h$. Let $h \rightarrow 0$. By continuity of f , it follows that the limit of the right side exists and so

$$\lim_{h \rightarrow 0} \frac{F(x+h) - F(x)}{h} = \lim_{h \rightarrow 0} f(x_h) = f(x)$$

If x is either end point, the argument is the same except you have to pay attention to the sign of h so that both x and $x+h$ are in $[a, b]$. Thus F is continuous on $[a, b]$ and F' exists on (a, b) so if G is an antiderivative,

$$\int_a^b f(t) dt \equiv F(b) - F(a) = G(b) - G(a)$$

The claim that the integral is linear is obvious from this. Indeed, if $F' = f, G' = g$,

$$\begin{aligned} \int_a^b (\alpha f(t) + \beta g(t)) dt &= \alpha F(b) + \beta G(b) - (\alpha F(a) + \beta G(a)) \\ &= \alpha (F(b) - F(a)) + \beta (G(b) - G(a)) \\ &= \alpha \int_a^b f(t) dt + \beta \int_a^b g(t) dt \end{aligned}$$

If $f \geq 0$, then the mean value theorem implies that for some

$$t \in (a, b), F(b) - F(a) = \int_a^b f dx = f(t)(b-a) \geq 0.$$

Thus

$$\int_a^b (|f| - f) dx \geq 0, \int_a^b (|f| + f) dx \geq 0$$

and so $\int_a^b |f| dx \geq \int_a^b f dx$ and $\int_a^b |f| dx \geq -\int_a^b f dx$ so this proves $\left| \int_a^b f dx \right| \leq \int_a^b |f| dx$. This, along with part 2 implies the other claim that $\left| \int_a^b f dx \right| \leq \int_a^b |f| dx$.

The last claim is obvious because an antiderivative of 1 is $F(x) = x$. ■

Note also that the usual change of variables theorem is available because if $F' = f$, then $f(g(x))g'(x) = \frac{d}{dx}F(g(x))$ so that, from the above proposition,

$$F(g(b)) - F(g(a)) = \int_{g(a)}^{g(b)} f(y) dy = \int_a^b f(g(x))g'(x) dx.$$

We usually let $y = g(x)$ and $dy = g'(x)dx$ and then change the limits as indicated above, equivalently we massage the expression to look like the above. Integration by parts also follows from differentiation rules.

Consider the iterated integral $\int_{a_1}^{b_1} \cdots \int_{a_p}^{b_p} \alpha x_1^{\alpha_1} \cdots x_p^{\alpha_p} dx_p \cdots dx_1$. It means just what it meant in calculus. You do the integral with respect to x_p first, keeping the other variables constant, obtaining a polynomial function of the other variables. Then you do this one with respect to x_{p-1} and so forth. Thus, doing the computation, it reduces to

$$\alpha \prod_{k=1}^p \left(\int_{a_k}^{b_k} x_k^{\alpha_k} dx_k \right) = \alpha \prod_{k=1}^p \left(\frac{b_k^{\alpha_k+1}}{\alpha_k+1} - \frac{a_k^{\alpha_k+1}}{\alpha_k+1} \right)$$

and the same thing would be obtained for any other order of the iterated integrals. Since each of these integrals is linear, it follows that if (i_1, \dots, i_p) is any permutation of $(1, \dots, p)$, then for any polynomial q ,

$$\int_{a_1}^{b_1} \cdots \int_{a_p}^{b_p} q(x_1, \dots, x_p) dx_p \cdots dx_1 = \int_{a_{i_1}}^{b_{i_1}} \cdots \int_{a_{i_p}}^{b_{i_p}} q(x_1, \dots, x_p) dx_{i_p} \cdots dx_{i_1}$$

Now let $f : \prod_{k=1}^p [a_k, b_k] \rightarrow \mathbb{R}$ be continuous. Then each iterated integral results in a continuous function of the remaining variables and so the iterated integral makes sense. For example, by Proposition 5.9.5, $\left| \int_c^d f(x, y) dy - \int_c^d f(\hat{x}, y) dy \right| =$

$$\left| \int_c^d (f(x, y) - f(\hat{x}, y)) dy \right| \leq \max_{y \in [c, d]} |f(x, y) - f(\hat{x}, y)| < \varepsilon$$

if $|x - \hat{x}|$ is sufficiently small, thanks to uniform continuity of f on the compact set $[a, b] \times [c, d]$. Thus it makes perfect sense to consider the iterated integral $\int_a^b \int_c^d f(x, y) dy dx$. Then using Proposition 5.9.5 on the iterated integrals along with Theorem 5.7.1, there exists a sequence of polynomials which converges to f uniformly $\{p_n\}$. Then applying Proposition 5.9.5 repeatedly,

$$\left| \int_{a_{i_p}}^{b_{i_1}} \cdots \int_{a_{i_p}}^{b_{i_p}} f(x) dx_p \cdots dx_1 - \int_{a_{i_p}}^{b_{i_1}} \cdots \int_{a_{i_p}}^{b_{i_p}} p_n(x) dx_p \cdots dx_1 \right|$$

$$\leq \|f - p_n\| \prod_{k=1}^p |b_k - a_k| \quad (5.5)$$

With this, it is easy to prove a rudimentary Fubini theorem valid for continuous functions.

Theorem 5.9.6 *$f : \prod_{k=1}^p [a_k, b_k] \rightarrow \mathbb{R}$ be continuous. Then for (i_1, \dots, i_p) any permutation of $(1, \dots, p)$,*

$$\int_{a_{i_p}}^{b_{i_1}} \cdots \int_{a_{i_p}}^{b_{i_p}} f(\mathbf{x}) dx_{i_p} \cdots dx_{i_1} = \int_{a_1}^{b_1} \cdots \int_{a_p}^{b_p} f(\mathbf{x}) dx_p \cdots dx_1$$

If $f \geq 0$, then the iterated integrals are nonnegative if each $a_k \leq b_k$.

Proof: Let $\|p_n - f\|_{\prod_{k=1}^p [a_k, b_k]} \rightarrow 0$ where p_n is a polynomial. Then from 5.5,

$$\begin{aligned} \int_{a_{i_1}}^{b_{i_1}} \cdots \int_{a_{i_p}}^{b_{i_p}} f(\mathbf{x}) dx_{i_p} \cdots dx_{i_1} &= \lim_{n \rightarrow \infty} \int_{a_{i_1}}^{b_{i_1}} \cdots \int_{a_{i_p}}^{b_{i_p}} p_n(\mathbf{x}) dx_{i_p} \cdots dx_{i_1} \\ &= \lim_{n \rightarrow \infty} \int_{a_1}^{b_1} \cdots \int_{a_p}^{b_p} p_n(\mathbf{x}) dx_p \cdots dx_1 = \int_{a_1}^{b_1} \cdots \int_{a_p}^{b_p} f(\mathbf{x}) dx_p \cdots dx_1 \blacksquare \end{aligned}$$

You could replace f with $f \mathcal{X}_G$ where $\mathcal{X}_G(\mathbf{x}) = 1$ if $\mathbf{x} \in G$ and 0 otherwise provided each section of G consisting of holding all variables constant but 1, consists of finitely many intervals. Thus you can integrate over all the usual sets encountered in beginning calculus.

5.10 The Stone Weierstrass Approximation Theorem

There is a profound generalization of the Weierstrass approximation theorem due to Stone. It has to be one of the most elegant things available. It holds on locally compact Hausdorff spaces but here I will show the version which is valid on compact sets. Later the more general version is discussed.

Definition 5.10.1 *\mathcal{A} is an algebra of functions if \mathcal{A} is a vector space and if whenever $f, g \in \mathcal{A}$ then $fg \in \mathcal{A}$.*

To begin with assume that the field of scalars is \mathbb{R} . This will be generalized later. Theorem 5.6.2 implies the following corollary. See Corollary 5.6.3.

Corollary 5.10.2 *The polynomials are dense in $C([a, b])$.*

Here is another approach to proving this theorem. It is the original approach used by Weierstrass. Let $m \in \mathbb{N}$ and consider c_m such that $\int_{-1}^1 c_m (1 - x^2)^m dx = 1$. Then

$$1 = 2 \int_0^1 c_m (1 - x^2)^m dx \geq 2c_m \int_0^1 (1 - x)^m dx = 2c_m \frac{1}{m+1}$$

so $c_m \leq m+1$. Then

$$\int_{\delta}^1 c_m (1 - x^2)^m dx + \int_{-1}^{-\delta} c_m (1 - x^2)^m dx \leq 2(m+1) (1 - \delta^2)^m$$

which converges to 0. Thus

$$\lim_{m \rightarrow \infty} \sup_{x \notin [-\delta, \delta]} c_m (1 - x^2)^m = 0 \quad (5.6)$$

Now let $\phi_n(t) \equiv c_m (1 - t^2)^m$. Consider $f \in C([-1, 1])$ and extend to let $f(x) = f(1)$ if $x > 1$ and $f(x) = f(-1)$ if $x < -1$ and define $p_m(x) \equiv \int_{-1}^1 f(x-t) \phi_m(t) dt$. Then

$$\begin{aligned} |p_m(x) - f(x)| &\leq \int_{-1}^1 |f(x-t) - f(x)| \phi_m(t) dt \leq \\ &\int_{-1}^1 \mathcal{X}_{[-\delta, \delta]}(t) |f(x-t) - f(x)| \phi_m(t) dt + \int_{-1}^1 \mathcal{X}_{[-1, 1] \setminus [-\delta, \delta]}(t) |f(x-t) - f(x)| \phi_m(t) dt \end{aligned}$$

Choose δ so small that if $|x-y| < \delta$, then $|f(x) - f(y)| < \varepsilon$. Also let $M \geq \max_x |f(x)|$. Then

$$\begin{aligned} |p_m(x) - f(x)| &\leq \varepsilon \int_{-1}^1 \phi_m(t) dt + 2M \int_{-1}^1 \mathcal{X}_{[-1, 1] \setminus [-\delta, \delta]}(t) \phi_m(t) dt \\ &= \varepsilon + 2M \int_{-1}^1 \mathcal{X}_{[-1, 1] \setminus [-\delta, \delta]}(t) \phi_m(t) dt \end{aligned}$$

From 5.6, The second term is no larger than $2M \int_{-1}^1 \mathcal{X}_{[-1, 1] \setminus [-\delta, \delta]}(t) \varepsilon dt \leq 4M\varepsilon$ whenever m is large enough. Hence, for large enough m , $\sup_{x \in [-1, 1]} |p_m(x) - f(x)| \leq (1 + 4M)\varepsilon$. Since ε is arbitrary, this shows that the functions p_m converge uniformly to f on $[-1, 1]$. However, p_m is actually a polynomial. To see this, change the variables and obtain

$$p_m(x) = \int_{x-1}^{x+1} f(t) \phi_m(x-t) dt$$

which will be a polynomial. To see this, note that a typical term is of the form

$$\int_{x-1}^{x+1} f(t) a(x-t)^k dt,$$

clearly a polynomial in x . This proves Corollary 5.10.2 in case $[a, b] = [-1, 1]$. In the general case, there is a linear one to one onto map $l : [-1, 1] \rightarrow [a, b]$.

$$l(t) = \frac{b-a}{2}(t+1) + a$$

Then if $f \in C([a, b])$, $f \circ l \in C([-1, 1])$. Hence there is a polynomial p such that

$$\max_{t \in [-1, 1]} |f \circ l(t) - p(t)| < \varepsilon$$

Then letting $t = l^{-1}(x) = \frac{2(x-a)}{b-a} - 1$, for $x \in [a, b]$, $\max_{x \in [a, b]} |f(x) - p(l^{-1}(x))| < \varepsilon$ but $x \rightarrow p(l^{-1}(x))$ is a polynomial. This gives an independent proof of that corollary. ■

The next result is the key to the profound generalization of the Weierstrass theorem due to Stone in which an interval will be replaced by a compact and later a locally compact set and polynomials will be replaced with elements of an algebra satisfying certain axioms.

Corollary 5.10.3 *On the interval $[-M, M]$, there exist polynomials p_n , $p_n(0) = 0$, and $\lim_{n \rightarrow \infty} \|p_n - |\cdot|\|_\infty = 0$. recall that $\|f\|_\infty \equiv \sup_{t \in [-M, M]} |f(t)|$.*

Proof: By Corollary 5.10.2 there exists a sequence of polynomials, $\{\tilde{p}_n\}$ such that $\tilde{p}_n \rightarrow |\cdot|$ uniformly. Then let $p_n(t) \equiv \tilde{p}_n(t) - \tilde{p}_n(0)$. ■

Definition 5.10.4 An algebra of functions, \mathcal{A} defined on A , annihilates no point of A if for all $x \in A$, there exists $g \in \mathcal{A}$ such that $g(x) \neq 0$. The algebra separates points if whenever $x_1 \neq x_2$, then there exists $g \in \mathcal{A}$ such that $g(x_1) \neq g(x_2)$.

The following generalization is known as the Stone Weierstrass approximation theorem.

Theorem 5.10.5 Let A be a compact topological space and let $\mathcal{A} \subseteq C(A; \mathbb{R})$ be an algebra of functions which separates points and annihilates no point. Then \mathcal{A} is dense in $C(A; \mathbb{R})$.

Proof: First here is a lemma.

Lemma 5.10.6 Let c_1 and c_2 be two real numbers and let $x_1 \neq x_2$ be two points of A . Then there exists a function $f_{x_1 x_2}$ such that

$$f_{x_1 x_2}(x_1) = c_1, f_{x_1 x_2}(x_2) = c_2.$$

Proof of the lemma: Let $g \in \mathcal{A}$ satisfy $g(x_1) \neq g(x_2)$. Such a g exists because the algebra separates points. Since the algebra annihilates no point, there exist functions h and k such that $h(x_1) \neq 0$, $k(x_2) \neq 0$. Then let $u \equiv gh - g(x_2)h$, $v \equiv gk - g(x_1)k$. It follows that $u(x_1) \neq 0$ and $u(x_2) = 0$ while $v(x_2) \neq 0$ and $v(x_1) = 0$. Let $f_{x_1 x_2} \equiv \frac{c_1 u}{u(x_1)} + \frac{c_2 v}{v(x_2)}$. This proves the lemma. Now continue the proof of Theorem 5.10.5.

First note that $\overline{\mathcal{A}}$ satisfies the same axioms as \mathcal{A} but in addition to these axioms, $\overline{\mathcal{A}}$ is closed. The closure of \mathcal{A} is taken with respect to the usual norm on $C(A)$,

$$\|f\|_\infty \equiv \max\{|f(x)| : x \in A\}.$$

Suppose $f \in \overline{\mathcal{A}}$ and suppose M is large enough that $\|f\|_\infty < M$. Using Corollary 5.10.3, let p_n be a sequence of polynomials such that

$$\|p_n - |\cdot|\|_\infty \rightarrow 0, p_n(0) = 0.$$

It follows that $p_n \circ f \in \overline{\mathcal{A}}$ and so $|f| \in \overline{\mathcal{A}}$ whenever $f \in \overline{\mathcal{A}}$. Also note that

$$\max(f, g) = \frac{|f - g| + (f + g)}{2}$$

$$\min(f, g) = \frac{(f + g) - |f - g|}{2}.$$

Therefore, this shows that if $f, g \in \overline{\mathcal{A}}$ then $\max(f, g), \min(f, g) \in \overline{\mathcal{A}}$. By induction, if $f_i, i = 1, 2, \dots, m$ are in $\overline{\mathcal{A}}$ then

$$\max(f_i, i = 1, 2, \dots, m), \min(f_i, i = 1, 2, \dots, m) \in \overline{\mathcal{A}}.$$

Now let $h \in C(A; \mathbb{R})$ and let $x \in A$. Use Lemma 5.10.6 to obtain f_{xy} , a function of $\overline{\mathcal{A}}$ which agrees with h at x and y . Letting $\varepsilon > 0$, there exists an open set $U(y)$ containing y such that

$$f_{xy}(z) > h(z) - \varepsilon \text{ if } z \in U(y).$$

Since A is compact, let $U(y_1), \dots, U(y_l)$ cover A . Let

$$f_x \equiv \max(f_{xy_1}, f_{xy_2}, \dots, f_{xy_l}).$$

Then $f_x \in \overline{\mathcal{A}}$ and $f_x(z) > h(z) - \varepsilon$ for all $z \in A$ and $f_x(x) = h(x)$. This implies that for each $x \in A$ there exists an open set $V(x)$ containing x such that for $z \in V(x)$, $f_x(z) < h(z) + \varepsilon$. Let $V(x_1), \dots, V(x_m)$ cover A and let $f \equiv \min(f_{x_1}, \dots, f_{x_m})$. Therefore, $f(z) < h(z) + \varepsilon$ for all $z \in A$ and since $f_x(z) > h(z) - \varepsilon$ for all $z \in A$, it follows $f(z) > h(z) - \varepsilon$ also and so $|f(z) - h(z)| < \varepsilon$ for all z . Since ε is arbitrary, this shows $h \in \overline{\mathcal{A}}$ and proves $\overline{\mathcal{A}} = C(A; \mathbb{R})$. ■

5.11 Connectedness in Normed Linear Space

The main result is that a ball in a normed linear space is connected. This is the next lemma. From this, it follows that for an open set, it is connected if and only if it is arcwise connected.

Lemma 5.11.1 *In a normed vector space, $B(z, r)$ is arcwise connected.*

Proof: This is easy from the convexity of the set. If $x, y \in B(z, r)$, then let $\gamma(t) = x + t(y - x)$ for $t \in [0, 1]$.

$$\begin{aligned} \|\mathbf{x} + t(\mathbf{y} - \mathbf{x}) - \mathbf{z}\| &= \|(1-t)(\mathbf{x} - \mathbf{z}) + t(\mathbf{y} - \mathbf{z})\| \\ &\leq (1-t)\|\mathbf{x} - \mathbf{z}\| + t\|\mathbf{y} - \mathbf{z}\| < (1-t)r + tr = r \end{aligned}$$

showing $\gamma(t)$ stays in $B(z, r)$. ■

Proposition 5.11.2 *If $X \neq \emptyset$ is arcwise connected, then it is connected.*

Proof: Let $p \in X$. Then by assumption, for any $x \in X$, there is an arc joining p and x . This arc is connected because it is the continuous image of an interval which is connected. Since x is arbitrary, every x is in a connected subset of X which contains p . Hence $C_p = X$ and so X is connected. ■

Theorem 5.11.3 *Let U be an open subset of a normed vector space. Then U is arcwise connected if and only if U is connected. Also the connected components of an open set are open sets.*

Proof: By Proposition 5.11.2 it is only necessary to verify that if U is connected and open in the context of this theorem, then U is arcwise connected. Pick $p \in U$. Say $x \in U$ satisfies \mathcal{P} if there exists a continuous function, $\gamma : [a, b] \rightarrow U$ such that $\gamma(a) = p$ and $\gamma(b) = x$.

$$A \equiv \{x \in U \text{ such that } x \text{ satisfies } \mathcal{P}\}$$

If $x \in A$, then Lemma 5.11.1 implies $B(x, r) \subseteq U$ is arcwise connected for small enough r . Thus letting $y \in B(x, r)$, there exist intervals, $[a, b]$ and $[c, d]$ and continuous functions having values in U , γ, η such that $\gamma(a) = p, \gamma(b) = x, \eta(c) = x$, and $\eta(d) = y$. Then let $\gamma_1 : [a, b + d - c] \rightarrow U$ be defined as

$$\gamma_1(t) \equiv \begin{cases} \gamma(t) & \text{if } t \in [a, b] \\ \eta(t + c - b) & \text{if } t \in [b, b + d - c] \end{cases}$$

Then it is clear that γ_1 is a continuous function mapping p to y and showing that $B(x, r) \subseteq A$. Therefore, A is open. $A \neq \emptyset$ because since U is open there is an open set, $B(p, \delta)$ containing p which is contained in U and is arcwise connected.

Now consider $B \equiv U \setminus A$. I claim this is also open. If B is not open, there exists a point $z \in B$ such that every open set containing z is not contained in B . Therefore, letting $B(z, \delta)$ be such that $z \in B(z, \delta) \subseteq U$, there exist points of A contained in $B(z, \delta)$. But then, a repeat of the above argument shows $z \in A$ also. Hence B is open and so if $B \neq \emptyset$, then $U = B \cup A$ and so U is separated by the two sets B and A contradicting the assumption that U is connected.

It remains to verify the connected components are open. Let $z \in C_p$ where C_p is the connected component determined by p . Then picking $B(z, \delta) \subseteq U$, $C_p \cup B(z, \delta)$ is connected and contained in U and so it must also be contained in C_p . Thus z is an interior point of C_p . ■

As an application, consider the following corollary.

Corollary 5.11.4 *Let $f : \Omega \rightarrow \mathbb{Z}$ be continuous where Ω is a connected open set in a normed vector space. Then f must be a constant.*

Proof: Suppose not. Then it achieves two different values, k and $l \neq k$. Then $\Omega = f^{-1}(l) \cup f^{-1}(\{m \in \mathbb{Z} : m \neq l\})$ and these are disjoint nonempty open sets which separate Ω . To see they are open, note

$$f^{-1}(\{m \in \mathbb{Z} : m \neq l\}) = f^{-1}\left(\bigcup_{m \neq l} \left(m - \frac{1}{6}, m + \frac{1}{6}\right)\right)$$

which is the inverse image of an open set while $f^{-1}(l) = f^{-1}\left((l - \frac{1}{6}, l + \frac{1}{6})\right)$ also an open set. ■

Definition 5.11.5 *An important concept in a vector space is the concept of convexity. A nonempty set K is called convex if whenever $x, y \in K$, it follows that for all $t \in [0, 1]$, $tx + (1-t)y \in K$ also. That is, the line segment joining the two points x, y is in K .*

5.12 Saddle Points*

A very useful idea in nonlinear analysis is the saddle point theorem also called the min max theorem. The proof of this theorem given here follows Brezis [8] which is where I found it. A real valued function f defined on a linear space is convex if

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

It is concave if the inequality is turned around. It can be shown that in finite dimensions, convex functions are automatically continuous, similar for concave functions. Recall the following definition of upper and lower semicontinuous functions defined on a metric space and having values in $[-\infty, \infty]$.

Definition 5.12.1 *A function is upper semicontinuous if whenever $x_n \rightarrow x$, it follows that $f(x) \geq \limsup_{n \rightarrow \infty} f(x_n)$ and it is lower semicontinuous if $f(x) \leq \liminf_{n \rightarrow \infty} f(x_n)$.*

The following lemma comes directly from the definition.

Lemma 5.12.2 *If \mathcal{F} is a set of functions which are upper semicontinuous, then $g(x) \equiv \inf \{f(x) : f \in \mathcal{F}\}$ is also upper semicontinuous. Similarly, if \mathcal{F} is a set of functions which are lower semicontinuous, then if $g(x) \equiv \sup \{f(x) : f \in \mathcal{F}\}$ it follows that g is lower semicontinuous.*

Note that in a metric space, the above definitions of upper and lower semicontinuity in terms of sequences are equivalent to the definitions that

$$f(x) \geq \limsup_{r \rightarrow 0} \{f(y) : y \in B(x, r)\}, \quad f(x) \leq \liminf_{r \rightarrow 0} \{f(y) : y \in B(x, r)\}$$

respectively.

Here is a technical lemma which will make the proof of the saddle point theorem shorter. It seems fairly interesting also.

Lemma 5.12.3 *Suppose $H : A \times B \rightarrow \mathbb{R}$ is strictly convex in the first argument and concave in the second argument where A, B are compact convex nonempty subsets of Banach spaces E, F respectively and $x \rightarrow H(x, y)$ is lower semicontinuous while $y \rightarrow H(x, y)$ is upper semicontinuous. Let*

$$H(g(y), y) \equiv \min_{x \in A} H(x, y)$$

Then $g(y)$ is uniquely defined and also for $t \in [0, 1]$,

$$\lim_{t \rightarrow 0} g(y + t(z - y)) = g(y).$$

Proof: First suppose both z, w yield the definition of $g(y)$. Then

$$H\left(\frac{z+w}{2}, y\right) < \frac{1}{2}H(z, y) + \frac{1}{2}H(w, y)$$

which contradicts the definition of $g(y)$. As to the existence of $g(y)$ this is nothing more than the theorem that a lower semicontinuous function defined on a compact set achieves its minimum.

Now consider the last claim about “hemicontinuity”, continuity along a line. For all $x \in A$, it follows from the definition of g that

$$H(g(y + t(z - y)), y + t(z - y)) \leq H(x, y + t(z - y))$$

By concavity of H in the second argument,

$$(1 - t)H(g(y + t(z - y)), y) + tH(g(y + t(z - y)), z) \quad (5.7)$$

$$\leq H(x, y + t(z - y)) \quad (5.8)$$

Now let $t_n \rightarrow 0$. Does $g(y + t_n(z - y)) \rightarrow g(y)$? Suppose not. By compactness, each of $g(y + t_n(z - y))$ is in a compact set and so there is a further subsequence, still denoted by t_n such that

$$g(y + t_n(z - y)) \rightarrow \hat{x} \in A$$

Then passing to a limit in 5.8, one obtains, using the upper semicontinuity in one and lower semicontinuity in the other the following inequality.

$$H(\hat{x}, y) \leq \liminf_{n \rightarrow \infty} (1 - t_n)H(g(y + t_n(z - y)), y) +$$

$$\begin{aligned}
& \liminf_{n \rightarrow \infty} t_n H(g(y + t_n(z - y)), z) \\
& \leq \liminf_{n \rightarrow \infty} \left(\begin{array}{c} (1 - t_n) H(g(y + t_n(z - y)), y) \\ + t_n H(g(y + t_n(z - y)), z) \end{array} \right) \\
& \leq \limsup_{n \rightarrow \infty} H(x, y + t_n(z - y)) \leq H(x, y)
\end{aligned}$$

This shows that $\hat{x} = g(y)$ because this holds for every x . Since $t_n \rightarrow 0$ was arbitrary, this shows that in fact

$$\lim_{t \rightarrow 0^+} g(y + t(z - y)) = g(y) \blacksquare$$

Now with this preparation, here is the min-max theorem.

Definition 5.12.4 *A norm is called strictly convex if whenever $x \neq y$,*

$$\left\| \frac{x + y}{2} \right\| < \frac{\|x\|}{2} + \frac{\|y\|}{2}$$

Theorem 5.12.5 *Let E, F be Banach spaces with E having a strictly convex norm. Also suppose that $A \subseteq E, B \subseteq F$ are compact and convex sets and that $H : A \times B \rightarrow \mathbb{R}$ is such that*

$$x \rightarrow H(x, y) \text{ is convex}$$

$$y \rightarrow H(x, y) \text{ is concave}$$

Assume that $x \rightarrow H(x, y)$ is lower semicontinuous and $y \rightarrow H(x, y)$ is upper semicontinuous. Then

$$\min_{x \in A} \max_{y \in B} H(x, y) = \max_{y \in B} \min_{x \in A} H(x, y)$$

This condition is equivalent to the existence of $(x_0, y_0) \in A \times B$ such that

$$H(x_0, y) \leq H(x_0, y_0) \leq H(x, y_0) \text{ for all } x, y \quad (5.9)$$

called a saddle point.

Proof: One part of the main equality is obvious.

$$\max_{y \in B} H(x, y) \geq H(x, y) \geq \min_{x \in A} H(x, y)$$

and so for each x ,

$$\max_{y \in B} H(x, y) \geq \max_{y \in B} \min_{x \in A} H(x, y)$$

and so

$$\min_{x \in A} \max_{y \in B} H(x, y) \geq \max_{y \in B} \min_{x \in A} H(x, y) \quad (5.10)$$

Next consider the other direction.

Define $H_\varepsilon(x, y) \equiv H(x, y) + \varepsilon \|x\|^2$ where $\varepsilon > 0$. Then H_ε is strictly convex in the first variable. This results from the observation that

$$\left\| \frac{x + y}{2} \right\|^2 < \left(\frac{\|x\| + \|y\|}{2} \right)^2 \leq \frac{1}{2} (\|x\|^2 + \|y\|^2),$$

Then by Lemma 5.12.3 there exists a unique $x \equiv g(y)$ such that

$$H_{\varepsilon}(g(y), y) \equiv \min_{x \in A} H_{\varepsilon}(x, y)$$

and also, whenever $y, z \in A$,

$$\lim_{t \rightarrow 0+} g(y + t(z - y)) = g(y).$$

Thus $H_{\varepsilon}(g(y), y) = \min_{x \in A} H_{\varepsilon}(x, y)$. But also this shows that $y \rightarrow H_{\varepsilon}(g(y), y)$ is the minimum of functions which are upper semicontinuous and so this function is also upper semicontinuous. Hence there exists y^* such that

$$\max_{y \in B} H_{\varepsilon}(g(y), y) = H_{\varepsilon}(g(y^*), y^*) = \max_{y \in B} \min_{x \in A} H_{\varepsilon}(x, y) \quad (5.11)$$

Thus from concavity in the second argument and what was just defined, for $t \in (0, 1)$,

$$\begin{aligned} H_{\varepsilon}(g(y^*), y^*) &\geq H_{\varepsilon}(g((1-t)y^* + ty), (1-t)y^* + ty) \\ &\geq (1-t)H_{\varepsilon}(g((1-t)y^* + ty), y^*) + tH_{\varepsilon}(g((1-t)y^* + ty), y) \\ &\geq (1-t)H_{\varepsilon}(g(y^*), y^*) + tH_{\varepsilon}(g((1-t)y^* + ty), y) \end{aligned} \quad (5.12)$$

This is because $\min_x H_{\varepsilon}(x, y^*) \equiv H_{\varepsilon}(g(y^*), y^*)$ so

$$H_{\varepsilon}(g((1-t)y^* + ty), y^*) \geq H_{\varepsilon}(g(y^*), y^*)$$

Then subtracting the first term on the right, one gets

$$tH_{\varepsilon}(g(y^*), y^*) \geq tH_{\varepsilon}(g((1-t)y^* + ty), y)$$

and cancelling the t ,

$$H_{\varepsilon}(g(y^*), y^*) \geq H_{\varepsilon}(g((1-t)y^* + ty), y)$$

Now apply Lemma 5.12.3 and let $t \rightarrow 0+$. This along with lower semicontinuity yields

$$H_{\varepsilon}(g(y^*), y^*) \geq \liminf_{t \rightarrow 0+} H_{\varepsilon}(g((1-t)y^* + ty), y) = H_{\varepsilon}(g(y^*), y) \quad (5.13)$$

Hence for every x, y

$$H_{\varepsilon}(x, y^*) \geq H_{\varepsilon}(g(y^*), y^*) \geq H_{\varepsilon}(g(y^*), y)$$

Thus

$$\min_x H_{\varepsilon}(x, y^*) \geq H_{\varepsilon}(g(y^*), y^*) \geq \max_y H_{\varepsilon}(g(y^*), y)$$

and so

$$\begin{aligned} \max_{y \in B} \min_{x \in A} H_{\varepsilon}(x, y) &\geq \min_x H_{\varepsilon}(x, y^*) \geq H_{\varepsilon}(g(y^*), y^*) \\ &\geq \max_y H_{\varepsilon}(g(y^*), y) \geq \min_{x \in A} \max_{y \in B} H_{\varepsilon}(x, y) \end{aligned}$$

Thus, letting $C \equiv \max \{\|x\| : x \in A\}$

$$\varepsilon C^2 + \max_{y \in B} \min_{x \in A} H(x, y) \geq \min_{x \in A} \max_{y \in B} H(x, y)$$

Since ε is arbitrary, it follows that

$$\max_{y \in B} \min_{x \in A} H(x, y) \geq \min_{x \in A} \max_{y \in B} H(x, y)$$

This proves the first part because it was shown above in 5.10 that

$$\min_{x \in A} \max_{y \in B} H(x, y) \geq \max_{y \in B} \min_{x \in A} H(x, y)$$

Now consider 5.9 about the existence of a “saddle point” given the equality of min max and max min. Let

$$\alpha = \max_{y \in B} \min_{x \in A} H(x, y) = \min_{x \in A} \max_{y \in B} H(x, y)$$

Then from

$$y \rightarrow \min_{x \in A} H(x, y) \text{ and } x \rightarrow \max_{y \in B} H(x, y)$$

being upper semicontinuous and lower semicontinuous respectively, there exist y_0 and x_0 such that

$$\alpha = \min_{x \in A} H(x, y_0) = \max_{y \in B} \min_{x \in A} H(x, y) \overset{\text{minimum of u.s.c.}}{=} \min_{x \in A} \max_{y \in B} H(x, y) \overset{\text{maximum of l.s.c.}}{=} \max_{y \in B} H(x_0, y)$$

Then

$$\alpha = \max_{y \in B} H(x_0, y) \geq H(x_0, y_0), \quad \alpha = \min_{x \in A} H(x, y_0) \leq H(x_0, y_0)$$

so in fact $\alpha = H(x_0, y_0)$ and from the above equalities,

$$\begin{aligned} H(x_0, y_0) &= \alpha = \min_{x \in A} H(x, y_0) \leq H(x, y_0) \\ H(x_0, y_0) &= \alpha = \max_{y \in B} H(x_0, y) \geq H(x_0, y) \end{aligned}$$

and so $H(x_0, y) \leq H(x_0, y_0) \leq H(x, y_0)$. Thus if the min max condition holds, then there exists a saddle point, namely (x_0, y_0) .

Finally suppose there is a saddle point (x_0, y_0) where

$$H(x_0, y) \leq H(x_0, y_0) \leq H(x, y_0)$$

Then

$$\min_{x \in A} \max_{y \in B} H(x, y) \leq \max_{y \in B} H(x_0, y) \leq H(x_0, y_0) \leq \min_{x \in A} H(x, y_0) \leq \max_{y \in B} \min_{x \in A} H(x, y)$$

However, as noted above, it is always the case that

$$\max_{y \in B} \min_{x \in A} H(x, y) \leq \min_{x \in A} \max_{y \in B} H(x, y) \quad \blacksquare$$

What was really needed? You needed compactness of A, B and these sets needed to be in a linear space. Of course there needed to be a norm for which $x \rightarrow \|x\|$ is strictly convex and lower semicontinuous, so the conditions given above are sufficient but maybe not necessary. You might try generalizing this much later after reading about weak topologies.

5.13 Exercises

1. Consider the metric space $C([0, T], \mathbb{R}^n)$ with the norm $\|\mathbf{f}\| \equiv \max_{x \in [0, T]} \|\mathbf{f}(x)\|_\infty$. Explain why the maximum exists. Show this is a complete metric space. **Hint:** If you have $\{\mathbf{f}_m\}$ a Cauchy sequence in $C([0, T], \mathbb{R}^n)$, then for each x , you have $\{\mathbf{f}_m(x)\}$ a Cauchy sequence in \mathbb{R}^n . Recall that this is a complete space. Thus there exists $\mathbf{f}(x) = \lim_{m \rightarrow \infty} \mathbf{f}_m(x)$. You must show that \mathbf{f} is continuous. This was in the section on the Ascoli Arzela theorem in more generality if you need an outline of how this goes. Write down the details for this case. Note how \mathbf{f} is in bold face. This means it is a function which has values in \mathbb{R}^n . $\mathbf{f}(t) = (f_1(t), f_2(t), \dots, f_n(t))$.
2. For $\mathbf{f} \in C([0, T], \mathbb{R}^n)$, you define the Riemann integral in the usual way using Riemann sums. Alternatively, you can define it as

$$\int_0^t \mathbf{f}(s) ds = \left(\int_0^t f_1(s) ds, \int_0^t f_2(s) ds, \dots, \int_0^t f_n(s) ds \right)$$

Then show that the following limit exists in \mathbb{R}^n for each $t \in (0, T)$.

$$\lim_{h \rightarrow 0} \frac{\int_0^{t+h} \mathbf{f}(s) ds - \int_0^t \mathbf{f}(s) ds}{h} = \mathbf{f}(t).$$

You should use the fundamental theorem of calculus from one variable calculus and the definition of the norm to verify this. As a review, in case we don't get to it in time, for \mathbf{f} defined on an interval $[0, T]$ and $s \in [0, T]$, $\lim_{t \rightarrow s} \mathbf{f}(t) = \mathbf{l}$ means that for all $\varepsilon > 0$, there exists $\delta > 0$ such that if $0 < |t - s| < \delta$, then $\|\mathbf{f}(t) - \mathbf{l}\|_\infty < \varepsilon$.

3. Suppose $f: \mathbb{R} \rightarrow \mathbb{R}$ and $f \geq 0$ on $[-1, 1]$ with $f(-1) = f(1) = 0$ and $f(x) < 0$ for all $x \notin [-1, 1]$. Can you use a modification of the proof of the Weierstrass approximation theorem for functions on an interval presented earlier to show that for all $\varepsilon > 0$ there exists a polynomial p , such that $|p(x) - f(x)| < \varepsilon$ for $x \in [-1, 1]$ and $p(x) \leq 0$ for all $x \notin [-1, 1]$?
4. A collection of functions \mathcal{F} of $C([0, T], \mathbb{R}^n)$ is said to be uniformly equicontinuous if for every $\varepsilon > 0$ there exists $\delta > 0$ such that if $\mathbf{f} \in \mathcal{F}$ and $|t - s| < \delta$, then $\|\mathbf{f}(t) - \mathbf{f}(s)\|_\infty < \varepsilon$. Thus the functions are uniformly continuous all at once. The single δ works for every pair t, s closer together than δ and for all functions $\mathbf{f} \in \mathcal{F}$. As an easy case, suppose there exists K such that for all $\mathbf{f} \in \mathcal{F}$, $\|\mathbf{f}(t) - \mathbf{f}(s)\|_\infty \leq K|t - s|$. Show that \mathcal{F} is uniformly equicontinuous. Now suppose \mathcal{G} is a collection of functions of $C([0, T], \mathbb{R}^n)$ which is bounded. That is, $\|\mathbf{f}\| = \max_{t \in [0, T]} \|\mathbf{f}(t)\|_\infty < M < \infty$ for all $\mathbf{f} \in \mathcal{G}$. Then let \mathcal{F} denote the functions which are of the form $\mathbf{F}(t) \equiv \mathbf{y}_0 + \int_0^t \mathbf{f}(s) ds$ where $\mathbf{f} \in \mathcal{G}$. Show that \mathcal{F} is uniformly equicontinuous. **Hint:** This is a really easy problem if you do the right things. Here is the way you should proceed. Remember the triangle inequality from one variable calculus which said that for $a < b$ $\left| \int_a^b f(s) ds \right| \leq \int_a^b |f(s)| ds$. Then $\left\| \int_a^b \mathbf{f}(s) ds \right\|_\infty = \max_i \left| \int_a^b f_i(s) ds \right| \leq \max_i \int_a^b |f_i(s)| ds \leq \int_a^b \|\mathbf{f}(s)\|_\infty ds$. Reduce to the case just considered using the assumption that these \mathbf{f} are bounded.
5. Suppose \mathcal{F} is a set of functions in $C([0, T], \mathbb{R}^n)$ which is uniformly bounded and uniformly equicontinuous as described above. Show it must be totally bounded.

6. †If $A \subseteq (X, d)$ is totally bounded, show that \bar{A} the closure of A is also totally bounded. In the above problem, explain why \mathcal{F} the closure of \mathcal{F} is compact. This uses the big theorem on compactness. Try and do this on your own, but if you get stuck, it is in the section on Arzela Ascoli theorem. When you have done this problem, you have proved the important part of the Arzela Ascoli theorem in the special case where the functions are defined on an interval. You can use this to prove one of the most important results in the theory of differential equations. This theorem is a really profound result because it gives compactness in a normed linear space which is **not finite dimensional**. Thus this is a non trivial generalization of the Heine Borel theorem.
7. Let $(X, \|\cdot\|)$ be a normed linear space. A set A is said to be **convex** if whenever $x, y \in A$ the line segment determined by these points given by $tx + (1-t)y$ for $t \in [0, 1]$ is also in A . Show that every open or closed ball is convex. Remember a closed ball is $D(x, r) \equiv \{\hat{x} : \|\hat{x} - x\| \leq r\}$ while the open ball is $B(x, r) \equiv \{\hat{x} : \|\hat{x} - x\| < r\}$. This should work just as easily in any normed linear space with any norm.
8. Let K be a nonempty closed and convex set in an inner product space $(X, |\cdot|)$ which is complete. For example, \mathbb{R}^n or any other finite dimensional inner product space. Let $y \notin K$ and let $\lambda = \inf\{|y - x| : x \in K\}$. Let $\{x_n\}$ be a minimizing sequence. That is $\lambda = \lim_{n \rightarrow \infty} |y - x_n|$. Explain why such a minimizing sequence exists. Next explain the following using the parallelogram identity in the above problem as follows.

$$\begin{aligned} \left| y - \frac{x_n + x_m}{2} \right|^2 &= \left| \frac{y}{2} - \frac{x_n}{2} + \frac{y}{2} - \frac{x_m}{2} \right|^2 \\ &= - \left| \frac{y}{2} - \frac{x_n}{2} - \left(\frac{y}{2} - \frac{x_m}{2} \right) \right|^2 + \frac{1}{2} |y - x_n|^2 + \frac{1}{2} |y - x_m|^2 \\ \text{Hence } \left| \frac{x_m - x_n}{2} \right|^2 &= - \left| y - \frac{x_n + x_m}{2} \right|^2 + \frac{1}{2} |y - x_n|^2 + \frac{1}{2} |y - x_m|^2 \\ &\leq -\lambda^2 + \frac{1}{2} |y - x_n|^2 + \frac{1}{2} |y - x_m|^2 \end{aligned}$$

Next explain why the right hand side converges to 0 as $m, n \rightarrow \infty$. Thus $\{x_n\}$ is a Cauchy sequence and converges to some $x \in X$. Explain why $x \in K$ and $|x - y| = \lambda$. Thus there exists a closest point in K to y . Next show that there is only one closest point. **Hint:** To do this, suppose there are two x_1, x_2 and consider $\frac{x_1 + x_2}{2}$ using the parallelogram law to show that this average works better than either of the two points which is a contradiction unless they are really the same point. This theorem is of enormous significance.

9. Let K be a closed convex nonempty set in a complete inner product space $(H, |\cdot|)$ (Hilbert space) and let $y \in H$. Denote the closest point to y by Px . Show that Px is characterized as being the solution to the following variational inequality

$$\operatorname{Re}(z - Py, y - Py) \leq 0$$

for all $z \in K$. That is, show that $x = Py$ if and only if $\operatorname{Re}(z - x, y - x) \leq 0$ for all $z \in K$. **Hint:** Let $x \in K$. Then, due to convexity, a generic thing in K is of the form $x + t(z - x), t \in [0, 1]$ for every $z \in K$. Then

$$|x + t(z - x) - y|^2 = |x - y|^2 + t^2 |z - x|^2 - t2\operatorname{Re}(z - x, y - x)$$

If $x = Px$, then the minimum value of this on the left occurs when $t = 0$. Function defined on $[0, 1]$ has its minimum at $t = 0$. What does it say about the derivative of this function at $t = 0$? Next consider the case that for some x the inequality $\operatorname{Re}(z - x, y - x) \leq 0$. Explain why this shows $x = Py$.

10. Using Problem 9 and Problem 8 show the projection map, P onto a closed convex subset is Lipschitz continuous with Lipschitz constant 1. That is $|Px - Py| \leq |x - y|$.
11. Suppose, in an inner product space, you know $\operatorname{Re}(x, y)$. Show that you also know $\operatorname{Im}(x, y)$. That is, give a formula for $\operatorname{Im}(x, y)$ in terms of $\operatorname{Re}(x, y)$. **Hint:**

$$(x, iy) = -i(x, y) = -i(\operatorname{Re}(x, y) + i\operatorname{Im}(x, y)) = -i\operatorname{Re}(x, y) + \operatorname{Im}(x, y)$$

while, by definition, $(x, iy) = \operatorname{Re}(x, iy) + i\operatorname{Im}(x, iy)$. Now consider matching real and imaginary parts.

12. Let $h > 0$ be given and let $\mathbf{f}(t, \mathbf{x}) \in \mathbb{R}^n$ for each $\mathbf{x} \in \mathbb{R}^n$. Also let $(t, \mathbf{x}) \rightarrow \mathbf{f}(t, \mathbf{x})$ be continuous and $\sup_{t, \mathbf{x}} \|\mathbf{f}(t, \mathbf{x})\|_\infty < C < \infty$. Let $\mathbf{x}_h(t)$ be a solution to the following

$$\mathbf{x}_h(t) = \mathbf{x}_0 + \int_0^t \mathbf{f}(s, \mathbf{x}_h(s-h)) ds$$

where $\mathbf{x}_h(s-h) \equiv \mathbf{x}_0$ if $s-h \leq 0$. Explain why there exists a solution. **Hint:** Consider the intervals $[0, h], [h, 2h]$ and so forth. Next explain why these functions $\{\mathbf{x}_h\}_{h>0}$ are equicontinuous and uniformly bounded. Now use the result of Problem 6 to argue that there exists a subsequence, still denoted by \mathbf{x}_h such that $\lim_{h \rightarrow 0} \mathbf{x}_h = \mathbf{x}$ in $C([0, T]; \mathbb{R}^n)$ as discussed in Problem 5. Use what you learned about the Riemann integral in single variable advanced calculus to explain why you can pass to a limit and conclude that $\mathbf{x}(t) = \mathbf{x}_0 + \int_0^t \mathbf{f}(s, \mathbf{x}(s)) ds$ **Hint:**

$$\begin{aligned} & \left\| \int_0^t \mathbf{f}(s, \mathbf{x}(s)) ds - \int_0^t \mathbf{f}(s, \mathbf{x}_h(s-h)) ds \right\|_\infty \\ & \leq \left\| \int_0^t \mathbf{f}(s, \mathbf{x}(s)) ds - \int_0^t \mathbf{f}(s, \mathbf{x}(s-h)) ds \right\|_\infty \\ & \quad + \left\| \int_0^t \mathbf{f}(s, \mathbf{x}(s-h)) ds - \int_0^t \mathbf{f}(s, \mathbf{x}_h(s-h)) ds \right\|_\infty \\ & \leq \int_0^T \|\mathbf{f}(s, \mathbf{x}(s)) - \mathbf{f}(s, \mathbf{x}(s-h))\|_\infty ds \\ & \quad + \int_0^T \|\mathbf{f}(s, \mathbf{x}(s-h)) - \mathbf{f}(s, \mathbf{x}_h(s-h))\|_\infty ds \end{aligned}$$

Now use Problem 2 to verify that $\mathbf{x}' = \mathbf{f}(t, \mathbf{x})$, $\mathbf{x}(0) = \mathbf{x}_0$. When you have done this, you will have proved the celebrated Peano existence theorem from ordinary differential equations.

13. Let $|\alpha| \equiv \sum_i \alpha_i$. Let \mathcal{G} denote all finite sums of functions of the form $p(\mathbf{x}) e^{-a|\mathbf{x}|^2}$ where $p(\mathbf{x})$ is a polynomial and $a > 0$. If you consider all real valued continuous functions defined on the closed ball $\overline{B(\mathbf{0}, R)}$ show that if f is such a function,

then for every $\varepsilon > 0$, there exists $g \in \mathcal{G}$ such that $\|f - g\|_\infty < \varepsilon$ where $\|h\|_\infty \equiv \max_{x \in \overline{B(0, R)}} |h(x)|$. Thus, from multi-variable calculus, every continuous function f is uniformly close to an infinitely differentiable function on any closed ball centered at 0.

14. Suppose now that $f \in C_0(\mathbb{R}^p)$. This means that f is everywhere continuous and that $\lim_{\|x\| \rightarrow \infty} |f(x)| = 0$. Show that for every $\varepsilon > 0$ there exists $g \in \mathcal{G}$ such that $\sup_{x \in \mathbb{R}^p} |f(x) - g(x)| < \varepsilon$. Thus you can approximate such a continuous function f uniformly on all of \mathbb{R}^p with a function which has infinitely many continuous partial derivatives. I assume the reader has had a beginning course in multi-variable calculus including partial derivatives. If not, a partial derivative is just a derivative with respect to one of the variables, fixing all the others.
15. In Problem 23 on Page 124, and $V \equiv \text{span}(f_{p_1}, \dots, f_{p_n}), f_r(x) \equiv x^r, x \in [0, 1]$ and $-\frac{1}{2} < p_1 < p_2 < \dots$ with $\lim_{k \rightarrow \infty} p_k = \infty$. The distance between f_m and V is

$$\frac{1}{\sqrt{2m+1}} \prod_{j \leq n} \frac{|m - p_j|}{(p_j + m + 1)} = d$$

Let $d_n = d$ so more functions are allowed to be included in V . Show that $\sum_n \frac{1}{p_n} = \infty$ if and only if $\lim_{n \rightarrow \infty} d_n = 0$. Explain, using the Weierstrass approximation theorem why this shows that if g is a function continuous on $[0, 1]$, then there is a function $\sum_{k=1}^N a_k f_{p_k}$ with $|g - \sum_{k=1}^N a_k f_{p_k}| < \varepsilon$. Here $|g|^2 \equiv \int_0^1 |g(x)|^2 dx$. This is Müntz's first theorem. **Hint:** $d_n \rightarrow 0$, if and only if $\ln d_n \rightarrow -\infty$ so you might want to arrange things so that this happens. You might want to use the fact that for $x \in [0, 1/2]$, $-x \geq \ln(1-x) \geq -2x$. See [10] which is where I read this. That product is $\prod_{j \leq n} \left(1 - \left(1 - \frac{|m - p_j|}{(p_j + m + 1)}\right)\right)$ and so \ln of this expression is

$$\sum_{j=1}^n \ln \left(1 - \left(1 - \frac{|m - p_j|}{(p_j + m + 1)}\right)\right)$$

which is in the interval

$$\left[-2 \sum_{j=1}^n \left(1 - \frac{|m - p_j|}{(p_j + m + 1)}\right), -\sum_{j=1}^n \left(1 - \frac{|m - p_j|}{(p_j + m + 1)}\right)\right]$$

and so $d_n \rightarrow 0$ if and only if $\sum_{j=1}^\infty \left(1 - \frac{|m - p_j|}{(p_j + m + 1)}\right) = \infty$. Since $p_n \rightarrow \infty$ it suffices to consider the convergence of $\sum_j \left(1 - \frac{p_j - m}{(p_j + m + 1)}\right) = \sum_j \left(\frac{2m+1}{(p_j + m + 1)}\right)$. Now recall theorems from calculus.

16. For $f \in C([a, b]; \mathbb{R})$, real valued continuous functions, let $|f| \equiv \left(\int_a^b |f(t)|^2\right)^{1/2} \equiv (f, f)^{1/2}$ where $(f, g) \equiv \int_a^b f(x)g(x)dx$. Recall the Cauchy Schwarz inequality $|(f, g)| \leq |f||g|$. Now suppose $\frac{1}{2} < p_1 < p_2 < \dots$ where $\lim_{k \rightarrow \infty} p_k = \infty$. Let $V_n = \text{span}(1, f_{p_1}, f_{p_2}, \dots, f_{p_n})$. For $\|\cdot\|$ the uniform approximation norm, show that for every $g \in C([0, 1])$, there exists there exists a sequence of functions, $f_n \in V_n$ such that

$\|g - f_n\| \rightarrow 0$. This is the second Müntz theorem. **Hint:** Show that you can approximate $x \rightarrow x^m$ uniformly. To do this, use the above Müntz to approximate mx^{m-1} with $\sum_k c_k x^{p_k-1}$ in the inner product norm. $\int_0^1 |mx^{m-1} - \sum_{k=1}^n c_k x^{p_k-1}|^2 dx \leq \epsilon^2$. Then $x^m - \sum_{k=1}^n \frac{c_k}{p_k} x^{p_k} = \int_0^x (mt^{m-1} - \sum_{k=1}^n c_k t^{p_k-1}) dt$. Then

$$\left| x^m - \sum_{k=1}^n \frac{c_k}{p_k} x^{p_k} \right| \leq \int_0^x \left| mt^{m-1} - \sum_{k=1}^n c_k t^{p_k-1} \right| dt \leq \int_0^1 1 \left| mt^{m-1} - \sum_{k=1}^n c_k t^{p_k-1} \right| dt$$

Now use the Cauchy Schwarz inequality on that last integral to obtain

$$\max_{x \in [0,1]} \left| x^m - \sum_{k=1}^n \frac{c_k}{p_k} x^{p_k} \right| \leq \epsilon.$$

In case $m = 0$, there is nothing to show because 1 is in V_n . Explain why the result follows from this and the Weierstrass approximation theorem.

Chapter 6

Fixed Point Theorems

This is on fixed point theorems which feature the Brouwer fixed point theorem. This next block of material is a discussion of simplices and triangulations used to prove the Brouwer fixed point theorem in an elementary way. It features the famous Sperner's lemma and is based on very elementary concepts from linear algebra in an essential way. However, it is pretty technical stuff. This elementary proof is harder than those which come from other approaches like integration theory or degree theory. These other shorter ways of obtaining the Brouwer fixed point theorem from analytical methods are presented later. If desired, this chapter could be placed after the easier to prove version of the Brouwer fixed point theorem, Theorem 11.6.8 on Page 329 after sufficient integration theory has been presented. I like the approach presented in this chapter which is based on simplices because it is elementary and contains a method for locating a fixed point. It seems philosophically wrong to make this theorem depend on integration theory.

6.1 Simplices and Triangulations

Definition 6.1.1 Define an n simplex, denoted by $[x_0, \dots, x_n]$, to be the convex hull of the $n+1$ points, $\{x_0, \dots, x_n\}$ where $\{x_i - x_0\}_{i=1}^n$ are linearly independent. Thus

$$[x_0, \dots, x_n] \equiv \left\{ \sum_{i=0}^n t_i x_i : \sum_{i=0}^n t_i = 1, t_i \geq 0 \right\}.$$

Note that $\{x_j - x_m\}_{j \neq m}$ are also independent. I will call the $\{t_i\}$ just described the coordinates of a point x .

To see the last claim, suppose $\sum_{j \neq m} c_j (x_j - x_m) = 0$. Then you would have

$$\begin{aligned} c_0 (x_0 - x_m) + \sum_{j \neq m, 0} c_j (x_j - x_m) &= 0 \\ &= c_0 (x_0 - x_m) + \sum_{j \neq m, 0} c_j (x_j - x_0) + \left(\sum_{j \neq m, 0} c_j \right) (x_0 - x_m) = 0 \\ &= \sum_{j \neq m, 0} c_j (x_j - x_0) + \left(\sum_{j \neq m} c_j \right) (x_0 - x_m) \end{aligned}$$

Then you get $\sum_{j \neq m} c_j = 0$ and each $c_j = 0$ for $j \neq m, 0$. Thus $c_0 = 0$ also because the sum is 0 and all other $c_j = 0$.

Since $\{x_i - x_0\}_{i=1}^n$ is an independent set, the t_i used to specify a point in the convex hull are uniquely determined. If two of them are $\sum_{i=0}^n t_i x_i = \sum_{i=0}^n s_i x_i$. Then $\sum_{i=0}^n t_i (x_i - x_0) = \sum_{i=0}^n s_i (x_i - x_0)$ so $t_i = s_i$ for $i \geq 1$ by independence. Since the s_i and t_i sum to 1, it follows that also $s_0 = t_0$. If $n \leq 2$, the simplex is a triangle, line segment, or point. If $n \leq 3$, it is a tetrahedron, triangle, line segment or point.

Definition 6.1.2 If S is an n simplex. Then it is triangulated if it is the union of smaller sub-simplices, the triangulation, such that if S_1, S_2 are two simplices in the triangulation, with

$$S_1 \equiv [z_0^1, \dots, z_m^1], S_2 \equiv [z_0^2, \dots, z_p^2]$$

then

$$S_1 \cap S_2 = [x_{k_0}, \dots, x_{k_r}]$$

where $[x_{k_0}, \dots, x_{k_r}]$ is in the triangulation and

$$\{x_{k_0}, \dots, x_{k_r}\} = \{z_0^1, \dots, z_m^1\} \cap \{z_0^2, \dots, z_p^2\}$$

or else the two simplices do not intersect.

The following proposition is geometrically fairly clear. It will be used without comment whenever needed in the following argument about triangulations.

Proposition 6.1.3 Say $[x_1, \dots, x_r], [\hat{x}_1, \dots, \hat{x}_r], [z_1, \dots, z_r]$ are all $r-1$ simplices and

$$[x_1, \dots, x_r], [\hat{x}_1, \dots, \hat{x}_r] \subseteq [z_1, \dots, z_r]$$

and $[z_1, \dots, z_r, b]$ is an $r+1$ simplex and

$$[y_1, \dots, y_s] = [x_1, \dots, x_r] \cap [\hat{x}_1, \dots, \hat{x}_r] \quad (6.1)$$

where

$$\{y_1, \dots, y_s\} = \{x_1, \dots, x_r\} \cap \{\hat{x}_1, \dots, \hat{x}_r\} \quad (6.2)$$

Then

$$[x_1, \dots, x_r, b] \cap [\hat{x}_1, \dots, \hat{x}_r, b] = [y_1, \dots, y_s, b] \quad (6.3)$$

Proof: If you have $\sum_{i=1}^s t_i y_i + t_{s+1} b$ in the right side, the t_i summing to 1 and nonnegative, then it is obviously in both of the two simplices on the left because of 6.2. Thus $[x_1, \dots, x_r, b] \cap [\hat{x}_1, \dots, \hat{x}_r, b] \supseteq [y_1, \dots, y_s, b]$.

Now suppose $x_k = \sum_{j=1}^r t_j^k z_j$, $\hat{x}_k = \sum_{j=1}^r \hat{t}_j^k z_j$, as usual, the scalars adding to 1 and nonnegative.

Consider something in both of the simplices on the left in 6.3. Is it in the right? The element on the left is of the form

$$\sum_{\alpha=1}^r s_\alpha x_\alpha + s_{r+1} b = \sum_{\alpha=1}^r \hat{s}_\alpha \hat{x}_\alpha + \hat{s}_{r+1} b$$

where the s_α , are nonnegative and sum to one, similarly for \hat{s}_α . Thus

$$\sum_{j=1}^r \sum_{\alpha=1}^r s_\alpha t_j^\alpha z_j + s_{r+1} b = \sum_{\alpha=1}^r \sum_{j=1}^r \hat{s}_\alpha \hat{t}_j^\alpha z_j + \hat{s}_{r+1} b \quad (6.4)$$

Now observe that

$$\sum_j \sum_\alpha s_\alpha t_j^\alpha + s_{r+1} = \sum_\alpha \sum_j s_\alpha t_j^\alpha + s_{r+1} = \sum_\alpha s_\alpha + s_{r+1} = 1.$$

A similar observation holds for the right side of 6.4. By uniqueness of the coordinates in an $r+1$ simplex, and assumption that $[z_1, \dots, z_r, b]$ is an $r+1$ simplex, $\hat{s}_{r+1} = s_{r+1}$ and so

$$\sum_{\alpha=1}^r \frac{s_\alpha}{1 - s_{r+1}} x_\alpha = \sum_{\alpha=1}^r \frac{\hat{s}_\alpha}{1 - s_{r+1}} \hat{x}_\alpha$$

where $\sum_{\alpha} \frac{s_{\alpha}}{1-s_{r+1}} = \sum_{\alpha} \frac{\hat{s}_{\alpha}}{1-s_{r+1}} = 1$, which would say that both sides are a single element of $[\mathbf{x}_1, \dots, \mathbf{x}_r] \cap [\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_r] = [\mathbf{y}_1, \dots, \mathbf{y}_s]$ and this shows both are equal to something of the form $\sum_{i=1}^s t_i \mathbf{y}_i$, $\sum_i t_i = 1, t_i \geq 0$. Therefore,

$$\sum_{\alpha=1}^r \frac{s_{\alpha}}{1-s_{r+1}} \mathbf{x}_{\alpha} = \sum_{i=1}^s t_i \mathbf{y}_i, \quad \sum_{\alpha=1}^r s_{\alpha} \mathbf{x}_{\alpha} = \sum_{i=1}^s (1-s_{r+1}) t_i \mathbf{y}_i$$

It follows that

$$\sum_{\alpha=1}^r s_{\alpha} \mathbf{x}_{\alpha} + s_{r+1} \mathbf{b} = \sum_{i=1}^s (1-s_{r+1}) t_i \mathbf{y}_i + s_{r+1} \mathbf{b} \in [\mathbf{y}_1, \dots, \mathbf{y}_s, \mathbf{b}]$$

which proves the other inclusion. ■

Next I will explain why any simplex can be triangulated in such a way that all sub-simplices have diameter less than ε .

This is obvious if $n \leq 2$. Supposing it to be true for $n-1$, is it also so for n ? The barycenter \mathbf{b} of a simplex $[\mathbf{x}_0, \dots, \mathbf{x}_n]$ is $\frac{1}{n+1} \sum_i \mathbf{x}_i$. This point is not in the convex hull of any of the faces, those simplices of the form $[\mathbf{x}_0, \dots, \hat{\mathbf{x}}_k, \dots, \mathbf{x}_n]$ where the hat indicates \mathbf{x}_k has been left out. Thus, placing \mathbf{b} in the k^{th} position, $[\mathbf{x}_0, \dots, \mathbf{b}, \dots, \mathbf{x}_n]$ is a n simplex also. First note that $[\mathbf{x}_0, \dots, \hat{\mathbf{x}}_k, \dots, \mathbf{x}_n]$ is an $n-1$ simplex. To be sure $[\mathbf{x}_0, \dots, \mathbf{b}, \dots, \mathbf{x}_n]$ is an n simplex, we need to check that certain vectors are linearly independent. If

$$\mathbf{0} = \sum_{j=1}^{k-1} c_j (\mathbf{x}_j - \mathbf{x}_0) + a_k \left(\frac{1}{n+1} \sum_{i=0}^n \mathbf{x}_i - \mathbf{x}_0 \right) + \sum_{j=k+1}^n d_j (\mathbf{x}_j - \mathbf{x}_0)$$

then does it follow that $a_k = 0 = c_j = d_j$?

$$\mathbf{0} = \sum_{j=1}^{k-1} c_j (\mathbf{x}_j - \mathbf{x}_0) + a_k \frac{1}{n+1} \left(\sum_{i=0}^n (\mathbf{x}_i - \mathbf{x}_0) \right) + \sum_{j=k+1}^n d_j (\mathbf{x}_j - \mathbf{x}_0)$$

$$\begin{aligned} \mathbf{0} &= \sum_{j=1}^{k-1} \left(c_j + \frac{a_k}{n+1} \right) (\mathbf{x}_j - \mathbf{x}_0) + a_k \frac{1}{n+1} (\mathbf{x}_k - \mathbf{x}_0) \\ &\quad + \sum_{j=k+1}^n \left(d_j + \frac{a_k}{n+1} \right) (\mathbf{x}_j - \mathbf{x}_0) \end{aligned}$$

Thus $\frac{a_k}{n+1} = 0$ and each $c_j + \frac{a_k}{n+1} = 0 = d_j + \frac{a_k}{n+1}$ so each c_j and d_j are also 0. Thus, this is also an n simplex.

Actually, a little more is needed. Suppose $[\mathbf{y}_0, \dots, \mathbf{y}_{n-1}]$ is an $n-1$ simplex such that $[\mathbf{y}_0, \dots, \mathbf{y}_{n-1}] \subseteq [\mathbf{x}_0, \dots, \hat{\mathbf{x}}_k, \dots, \mathbf{x}_n]$. Why is $[\mathbf{y}_0, \dots, \mathbf{y}_{n-1}, \mathbf{b}]$ an n simplex? We know the vectors $\{\mathbf{y}_j - \mathbf{y}_0\}_{j=1}^{n-1}$ are independent and that $\mathbf{y}_j = \sum_{i \neq k} t_i^j \mathbf{x}_i$ where $\sum_{i \neq k} t_i^j = 1$ with each being nonnegative. Suppose

$$\sum_{j=1}^{n-1} c_j (\mathbf{y}_j - \mathbf{y}_0) + c_n (\mathbf{b} - \mathbf{y}_0) = \mathbf{0} \quad (6.5)$$

If $c_n = 0$, then by assumption, each $c_j = 0$. The proof goes by assuming $c_n \neq 0$ and deriving a contradiction. Assume then that $c_n \neq 0$. Then you can divide by it and obtain modified

constants, still denoted as c_j such that

$$\mathbf{b} = \frac{1}{n+1} \sum_{i=0}^n \mathbf{x}_i = \mathbf{y}_0 + \sum_{j=1}^{n-1} c_j (\mathbf{y}_j - \mathbf{y}_0)$$

Thus

$$\begin{aligned} \frac{1}{n+1} \sum_{i=0}^n \sum_{s \neq k} t_s^0 (\mathbf{x}_i - \mathbf{x}_s) &= \sum_{j=1}^{n-1} c_j (\mathbf{y}_j - \mathbf{y}_0) = \sum_{j=1}^{n-1} c_j \left(\sum_{s \neq k} t_s^j \mathbf{x}_s - \sum_{s \neq k} t_s^0 \mathbf{x}_s \right) \\ &= \sum_{j=1}^{n-1} c_j \left(\sum_{s \neq k} t_s^j (\mathbf{x}_s - \mathbf{x}_0) - \sum_{s \neq k} t_s^0 (\mathbf{x}_s - \mathbf{x}_0) \right) \end{aligned}$$

Modify the term on the left and simplify on the right to get

$$\frac{1}{n+1} \sum_{i=0}^n \sum_{s \neq k} t_s^0 ((\mathbf{x}_i - \mathbf{x}_0) + (\mathbf{x}_0 - \mathbf{x}_s)) = \sum_{j=1}^{n-1} c_j \left(\sum_{s \neq k} (t_s^j - t_s^0) (\mathbf{x}_s - \mathbf{x}_0) \right)$$

Thus,

$$\begin{aligned} \frac{1}{n+1} \sum_{i=0}^n \left(\sum_{s \neq k} t_s^0 \right) (\mathbf{x}_i - \mathbf{x}_0) &= \frac{1}{n+1} \sum_{i=0}^n \sum_{s \neq k} t_s^0 (\mathbf{x}_s - \mathbf{x}_0) \\ &\quad + \sum_{j=1}^{n-1} c_j \left(\sum_{s \neq k} (t_s^j - t_s^0) (\mathbf{x}_s - \mathbf{x}_0) \right) \end{aligned}$$

Then, taking out the $i = k$ term on the left yields

$$\begin{aligned} \frac{1}{n+1} \left(\sum_{s \neq k} t_s^0 \right) (\mathbf{x}_k - \mathbf{x}_0) &= -\frac{1}{n+1} \sum_{i \neq k} \left(\sum_{s \neq k} t_s^0 \right) (\mathbf{x}_i - \mathbf{x}_0) \\ &\quad - \frac{1}{n+1} \sum_{i=0}^n \sum_{s \neq k} t_s^0 (\mathbf{x}_s - \mathbf{x}_0) + \sum_{j=1}^{n-1} c_j \left(\sum_{s \neq k} (t_s^j - t_s^0) (\mathbf{x}_s - \mathbf{x}_0) \right) \end{aligned}$$

That on the right is a linear combination of vectors $(\mathbf{x}_r - \mathbf{x}_0)$ for $r \neq k$ so by independence, $\sum_{r \neq k} t_r^0 = 0$. However, each $t_r^0 \geq 0$ and these sum to 1 so this is impossible. Hence $c_n = 0$ after all and so each $c_j = 0$. Thus $[\mathbf{y}_0, \dots, \mathbf{y}_{n-1}, \mathbf{b}]$ is an n simplex.

Now in general, if you have an n simplex $[\mathbf{x}_0, \dots, \mathbf{x}_n]$, its diameter is the maximum of $|\mathbf{x}_k - \mathbf{x}_l|$ for all $k \neq l$. Consider $|\mathbf{b} - \mathbf{x}_j|$. It equals

$$\left| \sum_{i=0}^n \frac{1}{n+1} (\mathbf{x}_i - \mathbf{x}_j) \right| = \left| \sum_{i \neq j} \frac{1}{n+1} (\mathbf{x}_i - \mathbf{x}_j) \right| \leq \frac{n}{n+1} \text{diam}(S).$$

Consider the k^{th} face of S which is the simplex $[\mathbf{x}_0, \dots, \hat{\mathbf{x}}_k, \dots, \mathbf{x}_n]$. By induction, it has a triangulation into simplices which each have diameter no more than $\frac{n}{n+1} \text{diam}(S)$. Let these $n-1$ simplices be denoted by $\{S_1^k, \dots, S_{m_k}^k\}$. Then the simplices $\{[S_i^k, \mathbf{b}]\}_{i=1, k=1}^{m_k, n+1}$ are a triangulation of S such that $\text{diam}([S_i^k, \mathbf{b}]) \leq \frac{n}{n+1} \text{diam}(S)$. Do for $[S_i^k, \mathbf{b}]$ what was just done for S obtaining a triangulation of S as the union of what is obtained such that each simplex has diameter no more than $(\frac{n}{n+1})^2 \text{diam}(S)$. Continuing this way shows the existence of the desired triangulation. You simply do the process k times where $(\frac{n}{n+1})^k \text{diam}(S) < \varepsilon$.

6.2 Labeling Vertices

Next is a way to label the vertices. Let p_0, \dots, p_n be the first $n+1$ prime numbers. All vertices of a simplex $S = [x_0, \dots, x_n]$ having $\{x_k - x_0\}_{k=1}^n$ independent will be labeled with one of these primes. In particular, the vertex x_k will be labeled as p_k if the simplex is $[x_0, \dots, x_n]$. The “value” of a simplex will be the product of its labels. Triangulate this S .

Consider a 1 simplex whose vertices are from the vertices of S , the original n simplex $[x_{k_1}, x_{k_2}]$, label x_{k_1} as p_{k_1} and x_{k_2} as p_{k_2} . Then label all other vertices of this triangulation which occur on $[x_{k_1}, x_{k_2}]$ either p_{k_1} or p_{k_2} . Note that by independence of $\{x_k - x_r\}_{k \neq r}$, this cannot introduce an inconsistency because the segment cannot contain any other vertex of S . Then obviously there will be an odd number of simplices in this triangulation having value $p_{k_1} p_{k_2}$, that is a p_{k_1} at one end and a p_{k_2} at the other. Next consider the 2 simplices $[x_{k_1}, x_{k_2}, x_{k_3}]$ where the x_{k_i} are from S . Label all vertices of the triangulation which lie on one of these 2 simplices which have not already been labeled as either p_{k_1}, p_{k_2} , or p_{k_3} . Continue this way. This labels all vertices of the triangulation of S which have at least one coordinate zero. For the vertices of the triangulation which have all coordinates positive, the interior points of S , label these at random from any of p_0, \dots, p_n . (Essentially, this is the same idea. The “interior” points are the new ones not already labeled.) The idea is to show that there is an odd number of n simplices with value $\prod_{i=0}^n p_i$ in the triangulation and more generally, for each m simplex $[x_{k_1}, \dots, x_{k_{m+1}}]$, $m \leq n$ with the x_{k_i} an original vertex from S , there are an odd number of m simplices of the triangulation contained in $[x_{k_1}, \dots, x_{k_{m+1}}]$, having value $p_{k_1} \cdots p_{k_{m+1}}$. It is clear that this is the case for all such 1 simplices. For convenience, call such simplices $[x_{k_1}, \dots, x_{k_{m+1}}]$ m dimensional faces of S . An m simplex which is a subspace of this one will have the “correct” value if its value is $p_{k_1} \cdots p_{k_{m+1}}$.

Suppose that the labeling has produced an odd number of simplices of the triangulation contained in each m dimensional face of S which have the correct value. Take such an m dimensional face $[x_{j_1}, \dots, x_{j_{k+1}}]$. Consider $\hat{S} \equiv$

$$[x_{j_1}, \dots, x_{j_{k+1}}, x_{j_{k+2}}]$$

Then by induction, there is an odd number of k simplices on the s^{th} face

$$[x_{j_1}, \dots, \hat{x}_{j_s}, \dots, x_{j_{k+2}}]$$

having value $\prod_{i \neq s} p_{j_i}$. In particular, the face $[x_{j_1}, \dots, x_{j_{k+1}}, \hat{x}_{j_{k+2}}]$ has an odd number of simplices with value $\prod_{i \leq k+1} p_{j_i}$.

No simplex in any other face of \hat{S} can have this value by uniqueness of prime factorization. Pick a simplex on the face $[x_{j_1}, \dots, x_{j_{k+1}}, \hat{x}_{j_{k+2}}]$ which has correct value $\prod_{i \leq k+1} p_{j_i}$ and cross this simplex into \hat{S} . Continue crossing simplices having value $\prod_{i \leq k+1} p_{j_i}$ which have not been crossed till the process ends. It must end because there are an odd number of these simplices having value $\prod_{i \leq k+1} p_{j_i}$. If the process leads to the outside of \hat{S} , then one can always enter it again because there are an odd number of simplices with value $\prod_{i \leq k+1} p_{j_i}$ available and you will have used up an even number. Note that in this process, if you have a simplex with one side labeled $\prod_{i \leq k+1} p_{j_i}$, there is either one way in or out of this simplex or two depending on whether the remaining vertex is labeled $p_{j_{k+2}}$. When the process ends, the value of the simplex must be $\prod_{i=1}^{k+2} p_{j_i}$ because it will have the additional label $p_{j_{k+2}}$. Otherwise, there would be another route out of this, through the other side labeled $\prod_{i \leq k+1} p_{j_i}$. This identifies a simplex in the triangulation with value $\prod_{i=1}^{k+2} p_{j_i}$.

Then repeat the process with $\prod_{i \leq k+1} p_{j_i}$ valued simplices on $[x_{j_1}, \dots, x_{j_{k+1}}, \hat{x}_{j_{k+2}}]$ which have not been crossed. Repeating the process, entering from the outside, cannot deliver a $\prod_{i=1}^{k+2} p_{j_i}$ valued simplex encountered earlier because of what was just noted. There is either one or two ways to cross the simplices. In other words, the process is one to one in selecting a $\prod_{i \leq k+1} p_{j_i}$ simplex from crossing such a simplex on the selected face of \hat{S} . Continue doing this, crossing a $\prod_{i \leq k+1} p_{j_i}$ simplex on the face of \hat{S} which has not been crossed previously. This identifies an odd number of simplices having value $\prod_{i=1}^{k+2} p_{j_i}$. These are the ones which are “accessible” from the outside using this process. If there are any which are not accessible from outside, applying the same process starting inside one of these, leads to exactly one other inaccessible simplex with value $\prod_{i=1}^{k+2} p_{j_i}$. Hence these inaccessible simplices occur in pairs and so there are an odd number of simplices in the triangulation having value $\prod_{i=1}^{k+2} p_{j_i}$. We refer to this procedure of labeling as Sperner’s lemma. The system of labeling is well defined thanks to the assumption that $\{x_k - x_0\}_{k=1}^n$ is independent which implies that $\{x_k - x_i\}_{k \neq i}$ is also linearly independent. Thus there can be no ambiguity in the labeling of vertices on any “face” the convex hull of some of the original vertices of S . The following is a description of the system of labeling the vertices.

Lemma 6.2.1 *Let $[x_0, \dots, x_n]$ be an n simplex with $\{x_k - x_0\}_{k=1}^n$ independent, and let the first $n+1$ primes be p_0, p_1, \dots, p_n . Label x_k as p_k and consider a triangulation of this simplex. Labeling the vertices of this triangulation which occur on $[x_{k_1}, \dots, x_{k_s}]$ with any of p_{k_1}, \dots, p_{k_s} , beginning with all 1 simplices $[x_{k_1}, x_{k_2}]$ and then 2 simplices and so forth, there are an odd number of simplices $[y_{k_1}, \dots, y_{k_s}]$ of the triangulation contained in $[x_{k_1}, \dots, x_{k_s}]$ which have value $p_{k_1} \dots p_{k_s}$. This for $s = 1, 2, \dots, n$.*

A combinatorial method

We now give a brief discussion of the system of labeling for Sperner’s lemma from the point of view of counting numbers of faces rather than obtaining them with an algorithm. Let p_0, \dots, p_n be the first $n+1$ prime numbers. All vertices of a simplex $S = [x_0, \dots, x_n]$ having $\{x_k - x_0\}_{k=1}^n$ independent will be labeled with one of these primes. In particular, the vertex x_k will be labeled as p_k . The value of a simplex will be the product of its labels. Triangulate this S . Consider a 1 simplex coming from the original simplex $[x_{k_1}, x_{k_2}]$, label one end as p_{k_1} and the other as p_{k_2} . Then label all other vertices of this triangulation which occur on $[x_{k_1}, x_{k_2}]$ either p_{k_1} or p_{k_2} . The assumption of linear independence assures that no **other** vertex of S can be in $[x_{k_1}, x_{k_2}]$ so there will be no inconsistency in the labeling. Then obviously there will be an odd number of simplices in this triangulation having value $p_{k_1} p_{k_2}$, that is a p_{k_1} at one end and a p_{k_2} at the other. Suppose that the labeling has been done for all vertices of the triangulation which are on $[x_{j_1}, \dots, x_{j_{k+1}}]$,

$$\{x_{j_1}, \dots, x_{j_{k+1}}\} \subseteq \{x_0, \dots, x_n\}$$

any k simplex for $k \leq n-1$, and there is an odd number of simplices from the triangulation having value equal to $\prod_{i=1}^{k+1} p_{j_i}$. Consider $\hat{S} \equiv [x_{j_1}, \dots, x_{j_{k+1}}, x_{j_{k+2}}]$. Then by induction, there is an odd number of k simplices on the s^{th} face

$$[x_{j_1}, \dots, \hat{x}_{j_s}, \dots, x_{j_{k+1}}]$$

having value $\prod_{i \neq s} p_{j_i}$. In particular the face $[x_{j_1}, \dots, x_{j_{k+1}}, \hat{x}_{j_{k+2}}]$ has an odd number of simplices with value $\prod_{i=1}^{k+1} p_{j_i} := \hat{P}_k$. We want to argue that some simplex in the triangulation which is contained in \hat{S} has value $\hat{P}_{k+1} := \prod_{i=1}^{k+2} p_{j_i}$. Let Q be the number of $k+1$

simplices from the triangulation contained in \hat{S} which have two faces with value \hat{P}_k (A $k+1$ simplex has either 1 or 2 \hat{P}_k faces.) and let R be the number of $k+1$ simplices from the triangulation contained in \hat{S} which have exactly one \hat{P}_k face. These are the ones we want because they have value \hat{P}_{k+1} . Thus the number of faces having value \hat{P}_k which is described here is $2Q+R$. All interior \hat{P}_k faces being counted twice by this number. Now we count the total number of \hat{P}_k faces another way. There are P of them on the face $[x_{j_1}, \dots, x_{j_{k+1}}, \hat{x}_{j_{k+2}}]$ and by induction, P is odd. Then there are O of them which are not on this face. These faces got counted twice. Therefore,

$$2Q+R = P+2O$$

and so, since P is odd, so is R . Thus there is an odd number of \hat{P}_{k+1} simplices in \hat{S} .

We refer to this procedure of labeling as Sperner's lemma. The system of labeling is well defined thanks to the assumption that $\{x_k - x_0\}_{k=1}^n$ is independent which implies that $\{x_k - x_i\}_{k \neq i}$ is also linearly independent. Thus there can be no ambiguity in the labeling of vertices on any "face", the convex hull of some of the original vertices of S . Sperner's lemma is now a consequence of this discussion.

6.3 The Brouwer Fixed Point Theorem

$S \equiv [x_0, \dots, x_n]$ is a simplex in \mathbb{R}^n . Assume $\{x_i - x_0\}_{i=1}^n$ are linearly independent. Thus a typical point of S is of the form $\sum_{i=0}^n t_i x_i$ where the t_i are uniquely determined and the map $x \rightarrow t$ is continuous from S to the compact set

$$\{t \in \mathbb{R}^{n+1} : \sum t_i = 1, t_i \geq 0\}$$

The map $t \rightarrow x$ is one to one and clearly continuous. Since S is compact, it follows that the inverse map is also continuous. This is a general consideration but what follows is a short explanation why this is so in this specific example.

To see this, suppose $x^k \rightarrow x$ in S . Let $x^k \equiv \sum_{i=0}^n t_i^k x_i$ with x defined similarly with t_i^k replaced with t_i , $x \equiv \sum_{i=0}^n t_i x_i$. Then

$$x^k - x_0 = \sum_{i=0}^n t_i^k x_i - \sum_{i=0}^n t_i^k x_0 = \sum_{i=1}^n t_i^k (x_i - x_0)$$

Thus

$$x^k - x_0 = \sum_{i=1}^n t_i^k (x_i - x_0), \quad x - x_0 = \sum_{i=1}^n t_i (x_i - x_0)$$

Say t_i^k fails to converge to t_i for all $i \geq 1$. Then there exists a subsequence, still denoted with superscript k such that for each $i = 1, \dots, n$, it follows that $t_i^k \rightarrow s_i$ where $s_i \geq 0$ and some $s_i \neq t_i$. But then, taking a limit, it follows that

$$x - x_0 = \sum_{i=1}^n s_i (x_i - x_0) = \sum_{i=1}^n t_i (x_i - x_0)$$

which contradicts independence of the $x_i - x_0$. It follows that for all $i \geq 1$, $t_i^k \rightarrow t_i$. Since they all sum to 1, this implies that also $t_0^k \rightarrow t_0$. Thus the claim about continuity is verified.

Let $f : S \rightarrow S$ be continuous. When doing f to a point x , one obtains another point of S denoted as $\sum_{i=0}^n s_i x_i$. Thus in this argument the scalars s_i will be the components after doing f to a point of S denoted as $\sum_{i=0}^n t_i x_i$.

Consider a triangulation of S such that all simplices in the triangulation have diameter less than ε . The vertices of the simplices in this triangulation will be labeled from p_0, \dots, p_n the first $n+1$ prime numbers. If $[y_0, \dots, y_n]$ is one of these simplices in the triangulation, each vertex is of the form $\sum_{l=0}^n t_l x_l$ where $t_l \geq 0$ and $\sum_l t_l = 1$. Let y_i be one of these vertices, $y_i = \sum_{l=0}^n t_l x_l$, the t_l being determined by y_i . Define $r_j \equiv s_j/t_j$ if $t_j > 0$ and ∞ if $t_j = 0$. Then $p(y_i)$ will be the label placed on y_i . To determine this label, let r_k be the smallest of these ratios. Then the label placed on y_i will be p_k where r_k is the smallest of all these extended nonnegative real numbers just described. If there is duplication, pick p_k where k is smallest. The value of the simplex will be the product of the labels. What does it mean for the value of the simplex to be $P_n \equiv p_0 p_1 \cdots p_n$? It means that each of the first $n+1$ primes is assigned to exactly one of the $n+1$ vertices of the simplex so each $r_j > 0$ and there are no repeats in the r_j .

Note that for the vertices which are on $[x_{i_1}, \dots, x_{i_m}]$, these will be labeled from the list $\{p_{i_1}, \dots, p_{i_m}\}$ because $t_k = 0$ for each of these and so $r_k = \infty$ unless $k \in \{i_1, \dots, i_m\}$. In particular, this scheme labels x_i as p_i .

By the Sperner's lemma procedure described above, there are an odd number of simplices having value $\prod_{i \neq k} p_i$ on the k^{th} face and an odd number of simplices in the triangulation of S for which the value of the simplex is $p_0 p_1 \cdots p_n \equiv P_n$. Thus if $[y_0, \dots, y_n]$ is one of these simplices, and $p(y_i)$ is the label for y_i , $\prod_{i=0}^n p(y_i) = \prod_{j=0}^n p_j \equiv P_n$.

What is r_k , the smallest of those ratios in determining a label? Could it be larger than 1? r_k is certainly finite because at least some $t_j \neq 0$ since they sum to 1. Thus, if $r_k > 1$, you would have $s_k > t_k$. The s_j sum to 1 and so some $s_j < t_j$ since otherwise, the sum of the t_j equalling 1 would require the sum of the s_j to be larger than 1. Hence r_k was not really the smallest after all and so $r_k \leq 1$. Hence $s_k \leq t_k$. Thus if the value of a simplex is P_n , then for each vertex of the simplex, the smallest ratio associated with it is of the form $s_j/t_j \leq 1$ and each j gets used exactly once.

Let $\mathcal{S} \equiv \{S_1, \dots, S_m\}$ denote those simplices whose value is P_n . In other words, if $\{y_0, \dots, y_n\}$ are the vertices of one of these simplices in \mathcal{S} , and $y_s = \sum_{i=0}^n t_i^s x_i$, $r_{k_s} \leq r_j$ for all $j \neq k_s$ and $\{k_0, \dots, k_n\} = \{0, \dots, n\}$. Let b denote the barycenter of $S_k = [y_0, \dots, y_n]$. $b \equiv \sum_{i=0}^n \frac{1}{n+1} y_i$

Do the same system of labeling for each n simplex in a sequence of triangulations where the diameters of the simplices in the k^{th} triangulation are no more than 2^{-k} . Thus each of these triangulations has a n simplex having diameter no more than 2^{-k} which has value P_n . Let b_k be the barycenter of one of these n simplices having value P_n . By compactness, there is a subsequence, still denoted with the index k such that $b_k \rightarrow x$. This x is a fixed point.

Consider this last claim. $x = \sum_{i=0}^n t_i x_i$ and after applying f , the result is $\sum_{i=0}^n s_i x_i$. Then b_k is the barycenter of some σ_k having diameter no more than 2^{-k} which has value P_n . Say σ_k is a simplex having vertices $\{y_0^k, \dots, y_n^k\}$ and the value of $[y_0^k, \dots, y_n^k]$ is P_n . Thus also $\lim_{k \rightarrow \infty} y_i^k = x$. Re ordering these vertices if necessary, we can assume that the label for y_i^k is p_i which implies that the smallest ratio r_k is when $k = i$ and as noted above, this ratio is no larger than 1. Thus for each $i = 0, \dots, n$,

$$\frac{s_i}{t_i} \leq 1, \quad s_i \leq t_i$$

the i^{th} coordinate of $f(y_i^k)$ with respect to the original vertices of S decreases and each i is represented for $i = \{0, 1, \dots, n\}$. As noted above, $y_i^k \rightarrow x$ and so the i^{th} coordinate of y_i^k, t_i^k must converge to t_i . Hence if the i^{th} coordinate of $f(y_i^k)$ is denoted by s_i^k , $s_i^k \leq t_i^k$. By continuity of f , it follows that $s_i^k \rightarrow s_i$. Thus the above inequality is preserved on taking

$k \rightarrow \infty$ and so $0 \leq s_i \leq t_i$ this for each i and these s_i, t_i pertain to the single point x . But these s_i add to 1 as do the t_i and so in fact, $s_i = t_i$ for each i and so $f(x) = x$. This proves the following theorem which is the Brouwer fixed point theorem.

Theorem 6.3.1 *Let S be a simplex $[x_0, \dots, x_n]$ such that $\{x_i - x_0\}_{i=1}^n$ are independent. Also let $f : S \rightarrow S$ be continuous. Then there exists $x \in S$ such that $f(x) = x$.*

Corollary 6.3.2 *Let K be a closed convex bounded subset of \mathbb{R}^n . Let $f : K \rightarrow K$ be continuous. Then there exists $x \in K$ such that $f(x) = x$.*

Proof: Let S be a large simplex containing K and let P be the projection map onto K . See Problem 10 on Page 152 for the necessary properties of this projection map. Consider $g(x) \equiv f(Px)$. Then g satisfies the necessary conditions for Theorem 6.3.1 and so there exists $x \in S$ such that $g(x) = x$. But this says $x \in K$ and so $g(x) = f(x)$. ■

Definition 6.3.3 *A set B has the fixed point property if whenever $f : B \rightarrow B$ for f continuous, it follows that f has a fixed point.*

The proof of this corollary is pretty significant. By a homework problem, a closed convex set is a retract of \mathbb{R}^n . This is what it means when you say there is this continuous projection map which maps onto the closed convex set but does not change any point in the closed convex set. When you have a set A which is a subset of a set B which has the property that continuous functions $f : B \rightarrow B$ have fixed points, and there is a continuous map P from B to A which leaves points of A unchanged, then it follows that A will have the same “fixed point property”. You can probably imagine all sorts of sets which are retracts of closed convex bounded sets. Also, if you have a compact set B which has the fixed point property and $h : B \rightarrow h(B)$ with h one to one and continuous, it will follow that h^{-1} is continuous and that $h(B)$ will also have the fixed point property. This is very easy to show. This will allow further extensions of this theorem. This says that the fixed point property is topological.

Several of the following theorems are generalizations of the Brouwer fixed point theorem.

6.4 The Schauder Fixed Point Theorem

First we give a proof of the Schauder fixed point theorem which is an infinite dimensional generalization of the Brouwer fixed point theorem. This is a theorem which lives in Banach space. Recall that one of these is a complete normed vector space. There is also a version of this theorem valid in locally convex topological vector spaces where the theorem is sometimes called Tychonoff’s theorem. In infinite dimensions, the closed unit ball fails to have the fixed point property. Thus something more is needed to get a fixed point.

We let K be a closed convex subset of X a Banach space and let

f be continuous, $f : K \rightarrow K$, and $\overline{f(K)}$ is compact.

Lemma 6.4.1 *For each $r > 0$ there exists a finite set of points*

$$\{y_1, \dots, y_n\} \subseteq \overline{f(K)}$$

and continuous functions ψ_i defined on $\overline{f(K)}$ such that for $x \in \overline{f(K)}$,

$$\sum_{i=1}^n \psi_i(x) = 1, \quad (6.6)$$

$$\psi_i(x) = 0 \text{ if } x \notin B(y_i, r), \quad \psi_i(x) > 0 \text{ if } x \in B(y_i, r).$$

If

$$f_r(x) \equiv \sum_{i=1}^n y_i \psi_i(f(x)), \quad (6.7)$$

then whenever $x \in K$,

$$\|f(x) - f_r(x)\| \leq r.$$

Proof: Using the compactness of $\overline{f(K)}$ which implies this set is totally bounded, there exists an r net

$$\{y_1, \dots, y_n\} \subseteq \overline{f(K)} \subseteq K$$

such that $\{B(y_i, r)\}_{i=1}^n$ covers $\overline{f(K)}$. Let

$$\phi_i(y) \equiv (r - \|y - y_i\|)^+$$

Thus $\phi_i(y) > 0$ if $y \in B(y_i, r)$ and $\phi_i(y) = 0$ if $y \notin B(y_i, r)$. For $x \in \overline{f(K)}$, let

$$\psi_i(x) \equiv \phi_i(x) \left(\sum_{j=1}^n \phi_j(x) \right)^{-1}.$$

Then 6.6 is satisfied. Indeed the denominator is not zero because x is in one of the $B(y_i, r)$. Thus it is obvious that the sum of these $\psi_i(f(x))$ equals 1 for $x \in K$. Now let f_r be given by 6.7 for $x \in K$. For such x ,

$$f(x) - f_r(x) = \sum_{i=1}^n (f(x) - y_i) \psi_i(f(x))$$

Thus

$$\begin{aligned} f(x) - f_r(x) &= \sum_{\{i: f(x) \in B(y_i, r)\}} (f(x) - y_i) \psi_i(f(x)) \\ &\quad + \sum_{\{i: f(x) \notin B(y_i, r)\}} (f(x) - y_i) \psi_i(f(x)) \\ &= \sum_{\{i: f(x) - y_i \in B(0, r)\}} (f(x) - y_i) \psi_i(f(x)) = \\ &\quad \sum_{\{i: f(x) - y_i \in B(0, r)\}} (f(x) - y_i) \psi_i(f(x)) + \sum_{\{i: f(x) \notin B(y_i, r)\}} 0 \psi_i(f(x)) \in B(0, r) \end{aligned}$$

because $0 \in B(0, r)$, $B(0, r)$ is convex, and 6.6. It is just a convex combination of things in $B(0, r)$. ■

Note that we could have had the y_i in $f(K)$ in addition to being in $\overline{f(K)}$. This would make it possible to eliminate the assumption that K is closed later on. All you really need is that K is convex.

We think of f_r as an approximation to f . In fact it is uniformly within r of f on K . The next lemma shows that this f_r has a fixed point. This is the main result and comes from the Brouwer fixed point theorem in \mathbb{R}^n . This will be an approximate fixed point.

Lemma 6.4.2 *For each $r > 0$, there exists $x_r \in \text{convex hull of } \overline{f(K)} \subseteq K$ such that*

$$f_r(x_r) = x_r, \quad \|f_r(x) - f(x)\| < r \text{ for all } x$$

Proof: If $f_r(x_r) = x_r$ and $x_r = \sum_{i=1}^n a_i y_i$ for $\sum_{i=1}^n a_i = 1$ and the y_i described in the above lemma, we need

$$f_r(x_r) \equiv \sum_{i=1}^n y_i \psi_i(f(x_r)) = \sum_{j=1}^n y_j \psi_j \left(f \left(\sum_{i=1}^n a_i y_i \right) \right) = \sum_{j=1}^n a_j y_j = x_r.$$

Also, if this is satisfied, then we have the desired approximate fixed point.

This will be satisfied if for each $j = 1, \dots, n$,

$$a_j = \psi_j \left(f \left(\sum_{i=1}^n a_i y_i \right) \right); \quad (6.8)$$

so, let

$$\Sigma_{n-1} \equiv \left\{ \mathbf{a} \in \mathbb{R}^n : \sum_{i=1}^n a_i = 1, a_i \geq 0 \right\}$$

and let $h : \Sigma_{n-1} \rightarrow \Sigma_{n-1}$ be given by

$$h(\mathbf{a})_j \equiv \psi_j \left(f \left(\sum_{i=1}^n a_i y_i \right) \right).$$

Since h is a continuous function of \mathbf{a} , the Brouwer fixed point theorem applies and there exists a fixed point for h which is a solution to 6.8. ■

The following is the Schauder fixed point theorem.

Theorem 6.4.3 *Let K be a closed and convex subset of X , a normed linear space. Let $f : K \rightarrow K$ be continuous and suppose $\overline{f(K)}$ is compact. Then f has a fixed point.*

Proof: Recall that $f(x_r) - f_r(x_r) \in B(0, r)$ and $f_r(x_r) = x_r$ with $x_r \in \text{convex hull of } \overline{f(K)} \subseteq K$.

There is a subsequence, still denoted with subscript r with $r \rightarrow 0$ such that $f(x_r) \rightarrow x \in \overline{f(K)}$. **Note that the fact that K is convex is what makes f defined at x_r . x_r is in the convex hull of $\overline{f(K)} \subseteq K$. This is where we use K convex.** Then since f_r is uniformly close to f , it follows that $f(x_r) = x_r \rightarrow x$ also. Therefore,

$$f(x) = \lim_{r \rightarrow 0} f(x_r) = \lim_{r \rightarrow 0} f_r(x_r) = \lim_{r \rightarrow 0} x_r = x. \quad \blacksquare$$

We usually have in mind the mapping defined on a Banach space. However, the completeness was never used. Thus the result holds in a normed linear space.

There is a nice corollary of this major theorem which is called the Schaefer fixed point theorem or the Leray Schauder alternative principle [22].

Theorem 6.4.4 *Let $f : X \rightarrow X$ be a compact map. Then either*

1. *There is a fixed point for f for all $t \in [0, 1]$ or*

2. For every $r > 0$, there exists a solution to $x = tf(x)$ for $t \in (0, 1)$ such that $\|x\| > r$.

Proof: Suppose there is $t_0 \in [0, 1]$ such that $t_0 f$ has no fixed point. Then $t_0 \neq 0$. If $t_0 = 0$, then $t_0 f$ obviously has a fixed point. Thus $t_0 \in (0, 1]$. Then let r_M be the radial retraction onto $B(0, M)$.

$$r_M f(x) = M \frac{f(x)}{\|f(x)\|}$$

By Schauder's theorem there exists $x \in \overline{B(0, M)}$ such that $t_0 r_M f(x) = x$. Then if $\|f(x)\| \leq M$, r_M has no effect and so $t_0 f(x) = x$ which is assumed not to take place. Hence $\|f(x)\| > M$ and so $\|r_M f(x)\| = M$ so $\|x\| = t_0 M$. Also $t_0 r_M f(x) = t_0 M \frac{f(x)}{\|f(x)\|} = x$ and so $x = \hat{t} f(x)$, $\hat{t} = t_0 \frac{M}{\|f(x)\|} < 1$. Since M is arbitrary, it follows that the solutions to $x = tf(x)$ for $t \in (0, 1)$ are unbounded. It was just shown that there is a solution to $x = \hat{t} f(x)$, $\hat{t} < 1$ such that $\|x\| = t_0 M$ where M is arbitrary. Thus the second of the two alternatives holds. ■

As an example of the usefulness of the Schauder fixed point theorem, consider the following application to the theory of ordinary differential equations. In the context of this theorem, $X = C([0, T]; \mathbb{R}^n)$, a Banach space with norm given by

$$\|x\| \equiv \max \{|x(t)| : t \in [0, T]\}.$$

I assume the reader knows about the Riemann integral in what follows and the elementary fundamental theorem of calculus. More general versions of these things are presented later in the book.

Theorem 6.4.5 Let $f : [0, T] \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ be continuous and suppose there exists $L > 0$ such that for all $\lambda \in (0, 1)$, if

$$x' = \lambda f(t, x), \quad x(0) = x_0 \quad (6.9)$$

for all $t \in [0, T]$, then $\|x\| < L$. Then there exists a solution to

$$x' = f(t, x), \quad x(0) = x_0 \quad (6.10)$$

for $t \in [0, T]$.

Proof: Let $F : X \rightarrow X$ where X described above.

$$Fy(t) \equiv \int_0^t f(s, y(s) + x_0) ds$$

Let B be a bounded set in X . Then $|f(s, y(s) + x_0)|$ is bounded for $s \in [0, T]$ if $y \in B$. Say $|f(s, y(s) + x_0)| \leq C_B$. Hence $F(B)$ is bounded in X . Also, for $y \in B$, $s < t$,

$$|Fy(t) - Fy(s)| \leq \left| \int_s^t f(s, y(s) + x_0) ds \right| \leq C_B |t - s|$$

and so $F(B)$ is pre-compact by the Ascoli Arzela theorem. By the Schaefer fixed point theorem, there are two alternatives. Either there are unbounded solutions y to

$$\lambda F(y) = y$$

for various $\lambda \in (0, 1)$ or for all $\lambda \in [0, 1]$, there is a fixed point for λF . In the first case, there would be unbounded \mathbf{y}_λ solving

$$\mathbf{y}_\lambda(t) = \lambda \int_0^t \mathbf{f}(s, \mathbf{y}_\lambda(s) + \mathbf{x}_0) ds$$

Then let $\mathbf{x}_\lambda(s) \equiv \mathbf{y}_\lambda(s) + \mathbf{x}_0$ and you get $\|\mathbf{x}_\lambda\|$ also unbounded for various $\lambda \in (0, 1)$. The above implies

$$\mathbf{x}_\lambda(t) - \mathbf{x}_0 = \lambda \int_0^t \mathbf{f}(s, \mathbf{x}_\lambda(s)) ds$$

so $\mathbf{x}'_\lambda = \lambda \mathbf{f}(t, \mathbf{x}_\lambda)$, $\mathbf{x}_\lambda(0) = \mathbf{x}_0$ and these would be unbounded for $\lambda \in (0, 1)$ contrary to the assumption that there exists an estimate for these valid for all $\lambda \in (0, 1)$. Hence the first alternative must hold and hence there is $\mathbf{y} \in X$ such that $F\mathbf{y} = \mathbf{y}$. Then letting $\mathbf{x}(s) \equiv \mathbf{y}(s) + \mathbf{x}_0$, it follows that

$$\mathbf{x}(t) - \mathbf{x}_0 = \int_0^t \mathbf{f}(s, \mathbf{x}(s)) ds$$

and so \mathbf{x} is a solution to the differential equation on $[0, T]$. ■

Note that existence of a solution to the differential equation is not assumed, only estimates of possible solutions. These estimates are called *a-priori* estimates. Also note this is a global existence theorem, not a local one for a solution defined on only a small interval.

6.5 The Kakutani Fixed Point Theorem

Definition 6.5.1 If $A : X \rightarrow \mathcal{P}(Y)$ is a set-valued map, define the graph of A by

$$G(A) \equiv \{(x, y) : y \in Ax\}.$$

Consider a map A which maps \mathbb{C}^p to $\mathcal{P}(\mathbb{C}^p)$ which satisfies

$$A\mathbf{x} \text{ is compact and convex.} \quad (6.11)$$

and also the condition that if O is open and $O \supseteq A\mathbf{x}$, then there exists $\delta > 0$ such that if

$$\mathbf{y} \in B(\mathbf{x}, \delta), \text{ then } A\mathbf{y} \subseteq O. \quad (6.12)$$

This last condition is sometimes referred to as **upper semicontinuity**. In words, A is upper semicontinuous and has values which are compact and convex. This is equivalent to saying that if $A\mathbf{x} \in O$ and $\mathbf{x}_n \rightarrow \mathbf{x}$, then for large enough n , it follows that $A\mathbf{x}_n \subseteq O$.

With this definition, here is a lemma which has to do with the situation when the graph is closed.

Lemma 6.5.2 Let A satisfy 6.12. Then AK is a subset of a compact set whenever K is compact. Also the graph of A is closed if $A\mathbf{x}$ is closed.

Proof: Let $\mathbf{x} \in K$. Then $A\mathbf{x}$ is compact and contained in some open set whose closure is compact, $U_{\mathbf{x}}$. By assumption 6.12 there exists an open set $V_{\mathbf{x}}$ containing \mathbf{x} such that if $\mathbf{y} \in V_{\mathbf{x}}$, then $A\mathbf{y} \subseteq U_{\mathbf{x}}$. Let $V_{\mathbf{x}_1}, \dots, V_{\mathbf{x}_m}$ cover K . Then $AK \subseteq \bigcup_{k=1}^m \overline{U_{\mathbf{x}_k}}$, a compact set.

To see the graph of A is closed when $A\mathbf{x}$ is closed, let $\mathbf{x}_k \rightarrow \mathbf{x}$, $\mathbf{y}_k \rightarrow \mathbf{y}$ where $\mathbf{y}_k \in A\mathbf{x}_k$. Then letting $O = A\mathbf{x} + B(0, r)$ it follows from 6.12 that $\mathbf{y}_k \in A\mathbf{x}_k \subseteq O$ for all k large enough. Therefore, $\mathbf{y} \in A\mathbf{x} + B(0, 2r)$ and since $r > 0$ is arbitrary and $A\mathbf{x}$ is closed it follows $\mathbf{y} \in A\mathbf{x}$. ■

Also, there is a general consideration relative to upper semicontinuous functions.

Lemma 6.5.3 *If f is upper semicontinuous on some set K and g is continuous and defined on $f(K)$, then $g \circ f$ is also upper semicontinuous.*

Proof: Let $x_n \rightarrow x$ in K . Let $U \supseteq g \circ f(x)$. Is $g \circ f(x_n) \in U$ for all n large enough? We have $f(x) \in g^{-1}(U)$, an open set. Therefore, if n is large enough, $f(x_n) \in g^{-1}(U)$. It follows that for large enough n , $g \circ f(x_n) \in U$ and so $g \circ f$ is upper semicontinuous on K . ■

The next theorem is an application of the Brouwer fixed point theorem. First define an p simplex, denoted by $[x_0, \dots, x_p]$, to be the convex hull of the $p+1$ points, $\{x_0, \dots, x_p\}$ where $\{x_i - x_0\}_{i=1}^p$ are independent. Thus

$$[x_0, \dots, x_p] \equiv \left\{ \sum_{i=1}^p t_i x_i : \sum_{i=1}^p t_i = 1, t_i \geq 0 \right\}.$$

If $p \leq 2$, the simplex is a triangle, line segment, or point. If $p \leq 3$, it is a tetrahedron, triangle, line segment or point. A collection of simplices is a tiling of \mathbb{R}^p if \mathbb{R}^p is contained in their union and if S_1, S_2 are two simplices in the tiling, with

$$S_j = [x_0^j, \dots, x_p^j],$$

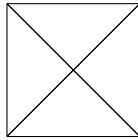
then

$$S_1 \cap S_2 = [x_{k_0}, \dots, x_{k_r}]$$

where

$$\{x_{k_0}, \dots, x_{k_r}\} \subseteq \{x_0^1, \dots, x_p^1\} \cap \{x_0^2, \dots, x_p^2\}$$

or else the two simplices do not intersect. The collection of simplices is said to be locally finite if, for every point, there exists a ball containing that point which also intersects only finitely many of the simplices in the collection. It is left to the reader to verify that for each $\varepsilon > 0$, there exists a locally finite tiling of \mathbb{R}^p which is composed of simplices which have diameters less than ε . The local finiteness ensures that for each ε the vertices have no limit point. To see how to do this, consider the case of \mathbb{R}^2 . Tile the plane with identical small squares and then form the triangles indicated in the following picture. It is clear something similar can be done in any dimension. Making the squares identical ensures that the little triangles are locally finite.



In general, you could consider $[0, 1]^p$. The point at the center is $(1/2, \dots, 1/2)$. Then there are $2p$ faces. Form the $2p$ pyramids having this point along with the 2^{p-1} vertices of the face. Then use induction on each of these faces to form smaller dimensional simplices tiling that face. Corresponding to each of these $2p$ pyramids, it is the union of the simplices whose vertices consist of the center point along with those of these new simplices tiling the chosen face. In general, you can write any p dimensional cube as the translate of a scaled $[0, 1]^p$. Thus one can express each of identical cubes as a tiling of $m(p)$ simplices of the

appropriate size and thereby obtain a tiling of \mathbb{R}^p with simplices. A ball will intersect only finitely many of the cubes and hence finitely many of the simplices. To get their diameters small as desired, just use $[0, r]^p$ instead of $[0, 1]^p$.

Thus one can give a function any value desired on these vertices and extend appropriately to the rest of the simplex and obtain a continuous function.

The Kakutani fixed point theorem is a generalization of the Brouwer fixed point theorem from continuous single valued maps to upper semicontinuous maps which have closed convex values.

Theorem 6.5.4 *Let K be a compact convex subset of \mathbb{R}^p and let $A : K \rightarrow \mathcal{P}(K)$ such that Ax is a closed convex subset of K and A is upper semicontinuous. Then there exists x such that $x \in Ax$. This is the “fixed point”.*

Proof: Let there be a locally finite tiling of \mathbb{R}^p consisting of simplices having diameter no more than ε . Let Px be the point in K which is closest to x . For each vertex x_k , pick $A_\varepsilon x_k \in APx_k$ and define A_ε on all of \mathbb{R}^p by the following rule. If

$$x \in [x_0, \dots, x_p],$$

so $x = \sum_{i=0}^p t_i x_i, t_i \in [0, 1], \sum_i t_i = 1$, then

$$A_\varepsilon x \equiv \sum_{k=0}^p t_k A_\varepsilon x_k.$$

Now by construction $A_\varepsilon x_k \in APx_k \in K$ and so A_ε is a continuous map defined on \mathbb{R}^p with values in K thanks to the local finiteness of the collection of simplices. By the Brouwer fixed point theorem A_ε has a fixed point x_ε in K , $A_\varepsilon x_\varepsilon = x_\varepsilon$.

$$x_\varepsilon = \sum_{k=0}^p t_k^\varepsilon A_\varepsilon x_k^\varepsilon, A_\varepsilon x_k^\varepsilon \in APx_k^\varepsilon \subseteq K$$

where a simplex containing x_ε is

$$[x_0^\varepsilon, \dots, x_p^\varepsilon], x_\varepsilon = \sum_{k=0}^p t_k^\varepsilon x_k^\varepsilon$$

Also, $x_\varepsilon \in K$ and is closer than ε to each x_k^ε so each x_k^ε is within ε of K . It follows that for each k , $|Px_k^\varepsilon - x_k^\varepsilon| < \varepsilon$ and so

$$\lim_{\varepsilon \rightarrow 0} |Px_k^\varepsilon - x_k^\varepsilon| = 0$$

By compactness of K , there exists a subsequence, still denoted with the subscript of ε such that for each k , the following convergences hold as $\varepsilon \rightarrow 0$

$$t_k^\varepsilon \rightarrow t_k, A_\varepsilon x_k^\varepsilon \rightarrow y_k, Px_k^\varepsilon \rightarrow z_k, x_k^\varepsilon \rightarrow z_k$$

Any pair of the x_k^ε are within ε of each other. Hence, any pair of the Px_k^ε are within ε of each other because P reduces distances. Therefore, in fact, z_k does not depend on k .

$$\lim_{\varepsilon \rightarrow 0} Px_k^\varepsilon = \lim_{\varepsilon \rightarrow 0} x_k^\varepsilon = z, \quad \lim_{\varepsilon \rightarrow 0} x_\varepsilon = \lim_{\varepsilon \rightarrow 0} \sum_{k=0}^p t_k^\varepsilon x_k^\varepsilon = \sum_{k=0}^p t_k z = z$$

By upper semicontinuity of A , for all ε small enough,

$$APx_k^\varepsilon \subseteq Az + B(0, r)$$

In particular, since $A_\varepsilon x_k^\varepsilon \in APx_k^\varepsilon$,

$$A_\varepsilon x_k^\varepsilon \in Az + B(0, r) \text{ for } \varepsilon \text{ small enough}$$

Since r is arbitrary and Az is closed, it follows $y_k \in Az$. It follows that since K is closed,

$$x_\varepsilon \rightarrow z = \sum_{k=0}^p t_k y_k, \quad t_k \geq 0, \quad \sum_{k=0}^p t_k = 1$$

Now by convexity of Az and the fact just shown that $y_k \in Az$,

$$z = \sum_{k=0}^p t_k y_k \in Az$$

and so $z \in Az$. This is the fixed point. ■

One can replace \mathbb{R}^p with \mathbb{C}^p in the above theorem because it is essentially \mathbb{R}^{2p} . Also the theorem holds with no change for any finite dimensional normed linear space since these are homeomorphic to \mathbb{R}^p or \mathbb{C}^p .

6.6 Ekeland's Variational Principle

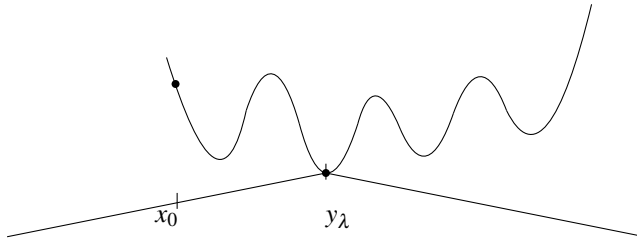
Recall the notation $X' = \mathcal{L}(X, \mathbb{R})$, the continuous linear functions mapping X to \mathbb{R} . This section deals with real Banach spaces. If you had complex ones, X' would denote $\mathcal{L}(X, \mathbb{C})$.

Definition 6.6.1 A function $\phi : X \rightarrow (-\infty, \infty]$ is called *proper* if it is not constantly equal to ∞ . Here X is assumed to be a complete metric space. The function ϕ is *lower semicontinuous* if

$$x_n \rightarrow x \text{ implies } \phi(x) \leq \liminf_{n \rightarrow \infty} \phi(x_n)$$

It is *bounded below* if there is some constant C such that $C \leq \phi(x)$ for all x .

The variational principle of Ekeland is the following theorem [22]. You start with an approximate minimizer x_0 . It says there is y_λ fairly close to x_0 such that if you subtract a “cone” from the value of ϕ at y_λ , then the resulting function is less than $\phi(x)$ for all $x \neq y_\lambda$.



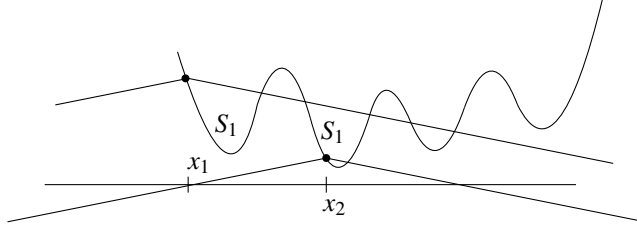
Theorem 6.6.2 Let X be a complete metric space and let $\phi : X \rightarrow (-\infty, \infty]$ be proper, lower semicontinuous and bounded below. Let x_0 be such that

$$\phi(x_0) \leq \inf_{x \in X} \phi(x) + \varepsilon$$

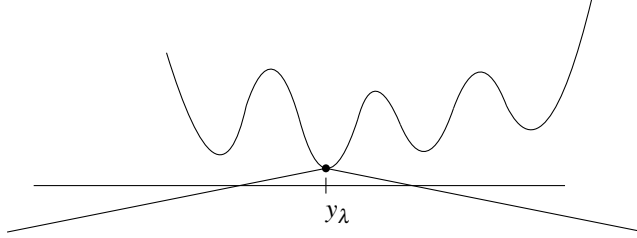
Then for every $\lambda > 0$ there exists a y_λ such that

1. $\phi(y_\lambda) \leq \phi(x_0)$
2. $d(y_\lambda, x_0) \leq \lambda$
3. $\phi(y_\lambda) - \frac{\varepsilon}{\lambda} d(x, y_\lambda) < \phi(x)$ for all $x \neq y_\lambda$

To motivate the proof, see the following picture which illustrates the first two steps. The S_i will be sets in X but are denoted symbolically by labeling them in $X \times (-\infty, \infty]$.



Then the end result of this iteration would be a picture like the following.



Thus you would have $\phi(y_\lambda) - \frac{\varepsilon}{\lambda} d(y_\lambda, x) \leq \phi(x)$ for all x which is seen to be what is wanted.

Proof: Let $x_1 = x_0$ and define $S_1 \equiv \{z \in X : \phi(z) \leq \phi(x_1) - \frac{\varepsilon}{\lambda} d(z, x_1)\}$. Then S_1 contains x_1 so it is nonempty. It is also clear that S_1 is a closed set. This follows from the lower semicontinuity of ϕ . Suppose

$$S_k \equiv \left\{ z \in X : \phi(z) \leq \phi(x_k) - \frac{\varepsilon}{\lambda} d(z, x_k) \right\}$$

where $x_k \in S_{k-1}$. Pick $x_{k+1} \in S_k$ and define S_{k+1} similarly. Will this yield a nested sequence of nonempty closed sets? Yes, it appears that it would because if $z \in S_k$ then

$$\begin{aligned} \phi(z) &\leq \phi(x_k) - \frac{\varepsilon}{\lambda} d(z, x_k) \leq \left(\phi(x_{k-1}) - \frac{\varepsilon}{\lambda} d(x_{k-1}, x_k) \right) - \frac{\varepsilon}{\lambda} d(z, x_k) \\ &\leq \phi(x_{k-1}) - \frac{\varepsilon}{\lambda} d(z, x_{k-1}) \end{aligned}$$

showing that z has what it takes to be in S_{k-1} . Thus we would obtain a sequence of nested, nonempty, closed sets according to this scheme.

Now here is how to choose the $x_k \in S_{k-1}$. Let $\phi(x_k) < \inf_{x \in S_{k-1}} \phi(x) + \frac{1}{2^k}$. Then for $z \in S_{n+1} \subseteq S_n$, $\phi(z) \leq \phi(x_{n+1}) - \frac{\varepsilon}{\lambda} d(z, x_{n+1})$ and so

$$\begin{aligned} \frac{\varepsilon}{\lambda} d(z, x_{n+1}) &\leq \phi(x_{n+1}) - \phi(z) \leq \inf_{x \in S_n} \phi(x) + \frac{1}{2^{n+1}} - \phi(z) \\ &\leq \phi(z) + \frac{1}{2^{n+1}} - \phi(z) = \frac{1}{2^{n+1}} \end{aligned}$$

Thus every $z \in S_{n+1}$ is within $\frac{1}{2^{n+1}}$ of the single point x_{n+1} and so the diameter of S_n converges to 0 as $n \rightarrow \infty$. By completeness of X , there exists a unique $y_\lambda \in \cap_n S_n$. Then it follows in particular that for $x_0 = x_1$ as above, $\phi(y_\lambda) \leq \phi(x_0) - \frac{\varepsilon}{\lambda} d(y_\lambda, x_0) \leq \phi(x_0)$ which verifies the first of the above conclusions.

As to the second, $\phi(x_0) \leq \inf_{x \in X} \phi(x) + \varepsilon$ and so, for any x ,

$$\phi(y_\lambda) \leq \phi(x_0) - \frac{\varepsilon}{\lambda} d(y_\lambda, x_0) \leq \phi(x) + \varepsilon - \frac{\varepsilon}{\lambda} d(y_\lambda, x_0),$$

this being true for $x = y_\lambda$. Hence $\frac{\varepsilon}{\lambda} d(y_\lambda, x_0) \leq \varepsilon$ and so $d(y_\lambda, x_0) \leq \lambda$.

Finally consider the third condition. If it does not hold, then there exists $z \neq y_\lambda$ such that $\phi(y_\lambda) \geq \phi(z) + \frac{\varepsilon}{\lambda} d(z, y_\lambda)$ so that $\phi(z) \leq \phi(y_\lambda) - \frac{\varepsilon}{\lambda} d(z, y_\lambda)$. But then, by the definition of y_λ as being in all the S_n , $\phi(y_\lambda) \leq \phi(x_n) - \frac{\varepsilon}{\lambda} d(x_n, y_\lambda)$ and so

$$\begin{aligned} \phi(z) &\leq \phi(x_n) - \frac{\varepsilon}{\lambda} (d(x_n, y_\lambda) + d(z, y_\lambda)) \\ &\leq \phi(x_n) - \frac{\varepsilon}{\lambda} d(x_n, z) \end{aligned}$$

Since n is arbitrary, this shows that $z \in \cap_n S_n$ but there is only one element of this intersection and it is y_λ so z must equal y_λ , a contradiction. ■

Note how if you make λ very small, you could pick ε very small such that the cone looks pretty flat. Of course, you can always consider an equivalent metric $\hat{d}(x, y) \equiv \frac{\varepsilon}{\lambda} d(x, y)$ in all of these considerations.

6.6.1 Cariste Fixed Point Theorem

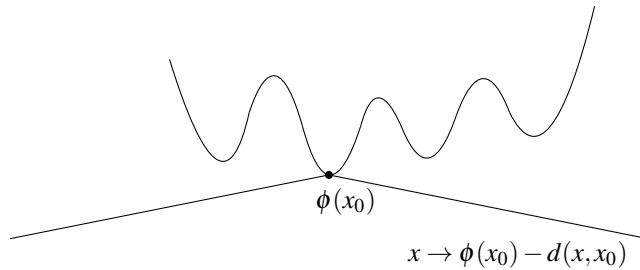
As mentioned in [22], the above result can be used to prove the Cariste fixed point theorem.

Theorem 6.6.3 *Let ϕ be lower semicontinuous, proper, and bounded below on a complete metric space X and let $F : X \rightarrow \mathcal{P}(X)$ be set valued such that $F(x) \neq \emptyset$ for all x . Also suppose that for each $x \in X$, there exists $y \in F(x)$ such that $\phi(y) \leq \phi(x) - d(x, y)$. Then there exists x_0 such that $x_0 \in F(x_0)$.*

Proof: In the above Ekeland variational principle, let $\varepsilon = 1 = \lambda$. Then there exists x_0 such that for all $y \neq x_0$

$$\phi(x_0) - d(y, x_0) < \phi(y), \text{ so } \phi(x_0) < \phi(y) + d(y, x_0) \quad (6.13)$$

for all $y \neq x_0$.



Suppose $x_0 \notin F(x_0)$. From the assumption, there is $y \in F(x_0)$ (so $y \neq x_0$) such that $\phi(y) \leq \phi(x_0) - d(x_0, y)$. Since $y \neq x_0$, it follows

$$\phi(y) + d(x_0, y) \leq \phi(x_0) < \phi(y) + d(y, x_0)$$

a contradiction. Hence $x_0 \in F(x_0)$ after all. ■

It is a funny theorem. It is easy to prove, but you look at it and wonder what it says. In fact, it implies the Banach fixed point theorem. If F is single valued, you would need to have a function ϕ such that for each x ,

$$\phi(F(x)) \leq \phi(x) - d(x, F(x))$$

and if you have such a ϕ then you can assert there is a fixed point for F . Suppose F is single valued and $d(Fx, Fy) \leq rd(x, y)$, $0 < r < 1$. Of course F has a fixed point using easier techniques. However, this also follows from this result. Let $\phi(x) = \frac{1}{1-r}d(x, F(x))$. Then is it true that for each x , there exists $y \in F(x)$ such that the inequality holds for all x ? Is

$$\frac{1}{1-r}d(F(x), F(F(x))) \leq \frac{1}{1-r}d(x, F(x)) - d(x, F(x))$$

Yes, this is certainly so because the right side reduces to $\frac{r}{1-r}d(x, F(x))$. Thus this fixed point theorem implies the usual Banach fixed point theorem.

The Ekeland variational principle says that when ϕ is lower semicontinuous proper and bounded below, there exists y such that

$$\phi(y) - d(x, y) < \phi(x) \text{ for all } x \neq y$$

In fact this can be proved from the Caristi fixed point theorem. Suppose the variational principle does not hold. This would mean that for all y there exists $x \neq y$ such that $\phi(y) - d(x, y) \geq \phi(x)$. Thus, for all x there exists $y \neq x$ such that $\phi(x) - d(x, y) \geq \phi(y)$. The inequality is preserved if $x = y$. Then let

$$F(x) \equiv \{y \neq x : \phi(x) - d(x, y) \geq \phi(y)\} \neq \emptyset$$

by assumption. This is the hypothesis for the Caristi fixed point theorem. Hence there exists x_0 such that $x_0 \in F(x_0) = \{y \neq x_0 : \phi(x_0) - d(x_0, y) \geq \phi(y)\}$ but this cannot happen because you can't have $x_0 \neq x_0$. Thus the Ekeland variational principle must hold after all.

6.6.2 A Density Result

There are several applications of the Ekeland variational principle. For more of them, see [22]. One of these is to show that there is a point where ϕ' is small assuming ϕ is bounded below, lower semicontinuous, and Gateaux differentiable, meaning that there exists $\phi'(x) \in X'$ such that if $v \in X$, then

$$\phi'(x)(v) \equiv \lim_{h \rightarrow 0} \frac{\phi(x + hv) - \phi(x)}{h}, \quad \phi'(x) \in X'$$

Here X is a real Banach space.

Theorem 6.6.4 *Let X be a Banach space and $\phi : X \rightarrow \mathbb{R}$ be Gateaux differentiable, bounded from below, and lower semicontinuous. Then for every $\varepsilon > 0$ there exists $x \in X$ such that*

$$\phi(x_\varepsilon) \leq \inf_{x \in X} \phi(x) + \varepsilon \text{ and } \|\phi'(x_\varepsilon)\|_{X'} \leq \varepsilon$$

Proof: From the Ekeland variational principle with $\lambda = 1$, there exists x_ε such that $\phi(x_\varepsilon) \leq \phi(x_0) \leq \inf_{x \in X} \phi(x) + \varepsilon$ and for all x , $\phi(x_\varepsilon) < \phi(x) + \varepsilon \|x - x_\varepsilon\|$. Then letting $x = x_\varepsilon + hv$ where $\|v\| = 1$, $\phi(x_\varepsilon + hv) - \phi(x_\varepsilon) > -\varepsilon|h|$. Let $h < 0$. Then divide by it to obtain $\frac{\phi(x_\varepsilon + hv) - \phi(x_\varepsilon)}{h} < \varepsilon$. Passing to a limit as $h \rightarrow 0$ yields $\phi'(x)(v) \leq \varepsilon$. Now v was arbitrary with norm 1 and so $\sup_{\|v\|=1} |\phi'(x_\varepsilon)(v)| = \|\phi'(x_\varepsilon)\| \leq \varepsilon$ ■

There is another very interesting application of the Ekeland variational principle [22].

Theorem 6.6.5 *Let X be a real Banach space and $\phi : X \rightarrow \mathbb{R}$ be Gateaux differentiable, bounded from below, and lower semicontinuous. Also suppose there exists $a, c > 0$ such that*

$$a\|x\| - c \leq \phi(x) \text{ for all } x \in X$$

Then $\{\phi'(x) : x \in X\}$ is dense in the ball of X' centered at 0 with radius a . Here $\phi'(x) \in X'$ and is determined by

$$\phi'(x)(v) \equiv \lim_{h \rightarrow 0} \frac{\phi(x + hv) - \phi(x)}{h}$$

Proof: Let $x^* \in X'$, $\|x^*\| \leq a$. Let $\psi(x) = \phi(x) - x^*(x)$. This is lower semicontinuous. It is also bounded from below because

$$\psi(x) \geq \phi(x) - a\|x\| \geq (a\|x\| - c) - a\|x\| = -c$$

It is also clearly Gateaux differentiable and lower semicontinuous because the piece added in is actually continuous. It is clear that the Gateaux derivative is just $\phi'(x) - x^*$. By Theorem 6.6.4, there exists x_ε such that $\|\phi'(x_\varepsilon) - x^*\| \leq \varepsilon$ ■

Thus this theorem says that if $\phi(x) \geq a\|x\| - c$ where ϕ has the nice properties of the theorem, it follows that $\phi'(x)$ is dense in $B(0, a)$ in the dual space X' . It follows that if for every a , there exists c such that $\phi(x) \geq a\|x\| - c$ for all $x \in X$ then $\{\phi'(x) : x \in X\}$ is dense in X' . This proves the following lemma.

Lemma 6.6.6 *Let X be a real Banach space and $\phi : X \rightarrow \mathbb{R}$ be Gateaux differentiable, bounded from below, and lower semicontinuous. Suppose for all $a > 0$ there exists a $c > 0$ such that $\phi(x) \geq a\|x\| - c$ for all x . Then $\{\phi'(x) : x \in X\}$ is dense in X' .*

If the above holds, then $\frac{\phi(x)}{\|x\|} \geq a - \frac{c}{\|x\|}$ and so, since a is arbitrary, it must be the case that

$$\lim_{\|x\| \rightarrow \infty} \frac{\phi(x)}{\|x\|} = \infty. \quad (6.14)$$

In fact, this is sufficient to conclude that for each $a > 0$ there is $c > 0$ such that $\phi(x) \geq a\|x\| - c$. If not, there would exist $a > 0$ such that $\phi(x_n) < a\|x_n\| - n$. Let $-L$ be a lower bound for $\phi(x)$. Then $-L + n \leq a\|x_n\|$ and so $\|x_n\| \rightarrow \infty$. Now it follows that

$$a \geq \frac{\phi(x_n)}{\|x_n\|} + \frac{n}{\|x_n\|} \geq \frac{\phi(x_n)}{\|x_n\|} \quad (6.15)$$

which is a contradiction to 6.14. This proves the following interesting density theorem.

Theorem 6.6.7 *Let X be a real Banach space and $\phi : X \rightarrow \mathbb{R}$ be Gateaux differentiable, bounded from below, and lower semicontinuous. Also suppose the coercivity condition*

$$\lim_{\|x\| \rightarrow \infty} \frac{\phi(x)}{\|x\|} = \infty$$

Then $\{\phi'(x) : x \in X\}$ is dense in X' . Here $\phi'(x) \in X'$ and is determined by

$$\phi'(x)(v) \equiv \lim_{h \rightarrow 0} \frac{\phi(x + hv) - \phi(x)}{h}$$

6.7 Exercises

1. It was shown that in a finite dimensional normed linear space that the compact sets are exactly those which are closed and bounded. Explain why every finite dimensional normed linear space is complete.
2. In any normed linear space, show that $\text{span}(x_1, \dots, x_n)$ is closed. That is, the span of any finite set of vectors is always a closed subspace. **Hint:** Suppose you let $V = \text{span}(x_1, \dots, x_n)$ and let $v^n \rightarrow v$ be convergent sequence of vectors in V . What does this say about the coordinate maps? Remember these are linear maps into \mathbb{F} and so they are continuous.
3. It was shown that in a finite dimensional normed linear space that the compact sets are exactly those which are closed and bounded. What if you have an infinite dimensional normed linear space X ? Show that the unit ball $D(0, r) \equiv \{x : \|x\| \leq 1\}$ is **NEVER** compact even though it is closed and bounded. **Hint:** Suppose you have $\{x_i\}_{i=1}^n$ where $\|x_i - x_j\| \geq \frac{1}{2}$. Let $y \notin \text{span}(x_1, \dots, x_n)$, a closed subspace. Such a y exists because X is not finite dimensional. Explain why $\text{dist}(y, \text{span}(x_1, \dots, x_n)) > 0$. This depends on $\text{span}(x_1, \dots, x_n)$ being closed. Let $z \in \text{span}(x_1, \dots, x_n)$ such that $\|y - z\| \leq 2\text{dist}(y, \text{span}(x_1, \dots, x_n))$. Let $x_{n+1} \equiv \frac{y-z}{\|y-z\|}$. Then consider the following:

$$\|x_{n+1} - x_k\| = \left\| \frac{y - (z + \|y - z\| x_k)}{\|y - z\|} \right\| \geq \frac{\|y - (z + \|y - z\| x_k)\|}{2\text{dist}(y, \text{span}(x_1, \dots, x_n))}$$

What of $(z + \|y - z\| x_k)$? Where is it? Isn't it in $\text{span}(x_1, \dots, x_n)$? Explain why this yields a sequence of points of X which are spaced at least $1/2$ apart even though they are all in the closed unit ball.

4. Find an example of two 2×2 matrices A, B such that $\|AB\| < \|A\| \|B\|$. This refers to the operator norm taken with respect to the usual norm on \mathbb{R}^2 . **Hint:** Maybe make it easy on yourself and consider diagonal matrices.
5. Now let $V = C([0, 1])$ and let $T : V \rightarrow V$ be given by $Tf(x) \equiv \int_0^x f(t) dt$. Show that T is continuous and linear. Here the norm is

$$\|f\| \equiv \max \{|f(x)| : x \in [0, 1]\}.$$

Can you find $\|T\|$ where this is the operator norm defined by analogy to what was given in the chapter?

6. Show that in any metric space (X, d) , if U is an open set and if $x \in U$, then there exists $r > 0$ such that the closure of $B(x, r)$, $\overline{B(x, r)} \subseteq U$. This says, in topological terms, that (X, d) is regular. Is it always the case in a metric space that $\overline{B(x, r)} = \{y : d(y, x) \leq r\} \equiv D(0, r)$? Prove or disprove. **Hint:** In fact, the answer to the last question is no.

7. Let (X, d) be a complete metric space. Let $\{U_n\}$ be a sequence of dense open sets. This means that $B(x, r) \cap U_n \neq \emptyset$ for every $x \in X$, and $r > 0$. You know that $\cap_n U_n$ is not necessarily open. Show that it is nevertheless, dense. **Hint:** Let $D = \cap_n U_n$. You need to show that $B(x, r) \cap D \neq \emptyset$. There is a point $p_1 \in U_1 \cap B(x, r)$. Then there exists $r_1 < 1/2$ such that $\overline{B(p_1, r_1)} \subseteq U_1 \cap B(x, r)$. From the above problem, you can adjust r_1 such that $\overline{B(p_1, r_1)} \subseteq U_1 \cap B(x, r)$. Next there exists $p_2 \in B(p_1, r_1) \cap U_2$. Let $r_2 < 1/2^2$ be such that $\overline{B(p_2, r_2)} \subseteq B(p_1, r_1) \cap U_2 \cap U_1$. Continue this way. You get a nested sequence of closed sets $\{B_k\}$ such that the diameter of B_k is no more than $1/2^{k-1}$, the k^{th} being contained in $B(p_{k-1}, r_{k-1}) \cap \cap_{i=1}^{k-1} U_i$. Explain why there is a unique point in the intersection of these closed sets which is in $B(x, r) \cap \cap_{k=1}^{\infty} U_k$. Then explain why this shows that D is dense.
8. The countable intersection of open sets is called a G_δ set. Show that the rational numbers \mathbb{Q} is **NOT** a G_δ set in \mathbb{R} . In fact, show that no countable dense set can be a G_δ set. Show that \mathbb{N} is a G_δ set. It is not dense.
9. You have a function $f : (X, d) \rightarrow (Y, \rho)$. Define

$$\omega_\delta f(x) \equiv \sup \{ \rho(f(z), f(y)) : z, y \in B(x, \delta) \}$$

Then explain why $\lim_{\delta \rightarrow 0} \omega_\delta f(x) \equiv \omega f(x)$ exists. Explain why a function is continuous at x if and only if $\omega f(x) = 0$. Next show that the set of all x where $\omega f(x) = 0$ is a G_δ set. **Hint:** $\omega f(x) = 0$ if and only if x is in something like this: $\cap_{n=1}^{\infty} \cup_{k=1}^{\infty} [\omega_{1/k} f(x) < \frac{1}{n}]$. Explain this. Then explain why $\cup_{k=1}^{\infty} [\omega_{1/k} f(x) < \frac{1}{n}]$ is an open set.

10. Prove or disprove.
- If A is compact, then $\mathbb{R}^n \setminus A$ is connected. You might consider the case $n > 1$ and the case $n = 1$ separately.
 - If A is connected in \mathbb{R}^n , then $\mathbb{R}^n \setminus A$ is also connected.
 - If A is connected in \mathbb{R}^n , then either A is open or A is closed.
 - $\mathbb{R}^n \setminus B(0, 1)$ is connected. Two cases to consider: $n = 1$ and $n > 1$.
11. If A is a connected set in \mathbb{R}^n , and A is not a single point, show that every point of A is a limit point of A .
12. Consider the Cantor set. This is obtained by starting with $[0, 1]$ deleting $(1/3, 2/3)$ and then taking the two closed intervals which result and deleting the middle open third of each of these and continuing this way. Let J_k denote the union of the 2^k closed intervals which result at the k^{th} step of the construction. The Cantor set is $J \equiv \cap_{k=1}^{\infty} J_k$. Explain why J is a nonempty compact subset of \mathbb{R} . Show that every point of J is a limit point of J . Also show there exists a mapping from J onto $[0, 1]$ even though the sum of the lengths of the deleted open intervals is 1. Show that the Cantor set has empty interior. If $x \in J$, consider the connected component of x . Show that this connected component is just x .
13. You have a complete metric space (X, d) and a mapping $T : X \rightarrow X$ which satisfies $d(Tx, Ty) \leq rd(x, y)$, $0 < r < 1$. Show x, Tx, T^2x, \dots converges to a point $z \in X$ such that $Tz = z$. Next suppose you only know $\frac{d(Tx, x)}{1-r} < R$ and that on

$B(x, R), d(Tx, Ty) \leq rd(x, y)$ where $r < 1$ as above. Show that then $z \in B(x, R)$ and that in fact each $T^k x \in B(x, R)$. Show also there is no more than one such fixed point z on $B(x, R)$.

14. In Theorem 5.7.1 it is assumed f has values in \mathbb{F} . Show there is no change if f has values in V , a normed vector space provided you redefine the definition of a polynomial to be something of the form $\sum_{|\alpha| \leq m} a_\alpha x^\alpha$ where $a_\alpha \in V$.
15. How would you generalize the conclusion of Corollary 5.8.8 to include the situation where f has values in a finite dimensional normed vector space?
16. If f and g are real valued functions which are continuous on some set, D , show that $\min(f, g), \max(f, g)$ are also continuous. Generalize this to any finite collection of continuous functions. **Hint:** Note $\max(f, g) = \frac{|f-g|+f+g}{2}$. Now recall the triangle inequality which can be used to show $|\cdot|$ is a continuous function.
17. Find an example of a sequence of continuous functions defined on \mathbb{R}^n such that each function is nonnegative and each function has a maximum value equal to 1 but the sequence of functions converges to 0 pointwise on $\mathbb{R}^n \setminus \{0\}$, that is, the set of vectors in \mathbb{R}^n excluding 0 .
18. An open subset U of \mathbb{R}^n is arcwise connected if and only if U is connected. Consider the usual Cartesian coordinates relative to axes x_1, \dots, x_n . A square curve is one consisting of a succession of straight line segments each of which is parallel to some coordinate axis. Show an open subset U of \mathbb{R}^n is connected if and only if every two points can be joined by a square curve.
19. Let $x \rightarrow h(x)$ be a bounded continuous function. Show f is continuous for $f(x) = \sum_{n=1}^{\infty} \frac{h(nx)}{n^2}$.
20. Let S be a any countable subset of \mathbb{R}^n . Show there exists a function, f defined on \mathbb{R}^n which is discontinuous at every point of S but continuous everywhere else. **Hint:** This is real easy if you do the right thing. It involves the Weierstrass M test.
21. If f is any continuous function defined on K a sequentially compact subset of \mathbb{R}^n , show there exists a series of the form $\sum_{k=1}^{\infty} p_k$, where each p_k is a polynomial, which converges uniformly to f on $[a, b]$. **Hint:** You should use the Weierstrass approximation theorem to obtain a sequence of polynomials. Then arrange it so the limit of this sequence is an infinite sum.
22. Let K be a sequentially compact set in a normed vector space V and let $f: V \rightarrow W$ be continuous where W is also a normed vector space. Show $f(K)$ is also sequentially compact.
23. If f is uniformly continuous, does it follow that $|f|$ is also uniformly continuous? If $|f|$ is uniformly continuous does it follow that f is uniformly continuous? Answer the same questions with “uniformly continuous” replaced with “continuous”. Explain why.

24. Suppose S, T are linear maps on some finite dimensional vector space, S^{-1} exists and let $\delta \in (0, 1)$. Then whenever $\|S - T\|$ is small enough, it follows that

$$\frac{|Tv|}{|Sv|} \in (1 - \delta, 1 + \delta) \quad (6.16)$$

for all $v \neq 0$. Similarly if T^{-1} exists and $\|S - T\|$ is small enough,

$$\frac{|Tv|}{|Sv|} \in (1 - \delta, 1 + \delta).$$

Hint: For the first part, consider the new norm $\|v\| \equiv |S^{-1}v|$. Use equivalence of norms and simple estimates to establish 6.16.

25. Let σ be an r simplex. Then $[\sigma, b]$ will consist of all $(1 - \lambda)\sigma + \lambda b$ where $\lambda \in [0, 1]$. If $\sigma = [x_1, \dots, x_r]$, show that $[\sigma, b] = [x_1, \dots, x_r, b]$. Now if $\sigma_1, \sigma_2 \subseteq \sigma$ where $[\sigma, b]$ is an $r + 1$ simplex and each σ_i is an r simplex, show that $[\sigma_1, b] \cap [\sigma_2, b] = [\sigma_2 \cap \sigma_1, b]$.
26. Let $A : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be continuous and let $f \in \mathbb{R}^n$. Also let (\cdot, \cdot) denote the standard inner product in \mathbb{R}^n . Letting K be a closed and bounded and convex set, show that there exists $x \in K$ such that for all $y \in K$, $(f - Ax, y - x) \leq 0$. **Hint:** Show that this is the same as saying $P(f - Ax + x) = x$ for some $x \in K$ where here P is the projection map discussed above in Problem 10 on Page 152. Now use the Brouwer fixed point theorem. This little observation is called Browder's lemma. It is a fundamental result in nonlinear analysis.
27. ↑ In the above problem, suppose that you have a coercivity result which is

$$\lim_{\|x\| \rightarrow \infty} \frac{(Ax, x)}{\|x\|} = \infty.$$

Show that if you have this, then you don't need to assume the convex closed set is bounded. In case $K = \mathbb{R}^n$, and this coercivity holds, show that A maps onto \mathbb{R}^n .

28. Let $f : X \rightarrow [-\infty, \infty]$ where X is a Banach space. This is said to be lower semicontinuous if whenever $x_n \rightarrow x$, it follows that $f(x) \leq \liminf_{n \rightarrow \infty} f(x_n)$. Show that this is the same as saying that the epigraph of f is closed. Here we can make $X \times [-\infty, \infty]$ into a metric space in a natural way by using the product topology where the distance on $[-\infty, \infty]$ will be $d(\sigma, \alpha) \equiv |\arctan(\sigma) - \arctan(\alpha)|$. Here $\text{epi}(f) \equiv \{(x, \alpha) : \alpha \geq f(x)\}$. The function is upper semicontinuous if $\limsup_{n \rightarrow \infty} f(x_n) \leq f(x)$. What is a condition for f to be upper semicontinuous? Do you need a Banach space to do this? Would it be sufficient to let X be a metric space?
29. Explain why the supremum of lower semicontinuous functions is lower semicontinuous and the infimum of upper semicontinuous functions is upper semicontinuous.
30. Let K be a nonempty closed and convex subset of \mathbb{R}^n . Recall K is convex means that if $x, y \in K$, then for all $t \in [0, 1]$, $tx + (1 - t)y \in K$. Show that if $x \in \mathbb{R}^n$ there exists a unique $z \in K$ such that $|x - z| = \min\{|x - y| : y \in K\}$. This z will be denoted as Px . **Hint:** First note you do not know K is compact. Establish the parallelogram

identity if you have not already done so, $|u - v|^2 + |u + v|^2 = 2|u|^2 + 2|v|^2$. Then let $\{z_k\}$ be a minimizing sequence,

$$\lim_{k \rightarrow \infty} |z_k - x|^2 = \inf \{|x - y| : y \in K\} \equiv \lambda.$$

Now using convexity, explain why

$$\left| \frac{z_k - z_m}{2} \right|^2 + \left| x - \frac{z_k + z_m}{2} \right|^2 = 2 \left| \frac{x - z_k}{2} \right|^2 + 2 \left| \frac{x - z_m}{2} \right|^2$$

and then use this to argue $\{z_k\}$ is a Cauchy sequence. Then if z_i works for $i = 1, 2$, consider $(z_1 + z_2)/2$ to get a contradiction.

31. In Problem 30 show that Px satisfies and is in fact characterized as the solution to the following variational inequality. $(x - Px, y - Px) \leq 0$ for all $y \in K$. Then show that $|Px_1 - Px_2| \leq |x_1 - x_2|$. **Hint:** For the first part note that if $y \in K$, the function $t \rightarrow |x - (Px + t(y - Px))|^2$ achieves its minimum on $[0, 1]$ at $t = 0$. For the second part,

$$(x_1 - Px_1) \cdot (Px_2 - Px_1) \leq 0, (x_2 - Px_2) \cdot (Px_1 - Px_2) \leq 0.$$

Explain why $(x_2 - Px_2 - (x_1 - Px_1)) \cdot (Px_2 - Px_1) \geq 0$ and then use a some manipulations and the Cauchy Schwarz inequality to get the desired inequality. Thus P is called a retraction onto K .

32. Browder's lemma says: Let K be a convex closed and bounded set in \mathbb{R}^n and let $A : K \rightarrow \mathbb{R}^n$ be continuous and $f \in \mathbb{R}^n$. Then there exists $x \in K$ such that for all $y \in K$,

$$(f - Ax, y - x) \leq 0$$

show this is true. **Hint:** Consider $x \rightarrow P(f - Ax + x)$ where P is the projection onto K . If there is a fixed point of this mapping, then $P(f - Ax + x) = x$. Now consider the variational inequality satisfied. This little lemma is the basis for a whole lot of nonlinear analysis involving nonlinear operators of various kinds.

33. Generalize the above problem as follows. Let K be a convex closed and bounded set in \mathbb{R}^n and let $A : K \rightarrow \mathcal{P}(\mathbb{R}^n)$ be upper semi-continuous having closed bounded convex values and $f \in \mathbb{R}^n$. Then there exists $x \in K$ and $z \in Ax$ such that for all $y \in K$, $(f - z, y - x) \leq 0$ show this is true. Also show that if K is a closed convex and bounded set in E a finite dimensional normed linear space and $A : K \rightarrow \mathcal{P}(E')$ is upper semicontinuous having closed bounded convex values and $f \in E'$, then there exists $x \in K$ and $z \in Ax$ such that for all $y \in K$, $\langle f - z, y - x \rangle \leq 0$. **Hint:** Use the construction for the proof of the Kakutani fixed point theorem and the above Browder's lemma.
34. This problem establishes a remarkable result about existence for a system of inequalities based on the min max theorem, Theorem 5.12.5. Let E be a finite dimensional Banach space and let K be a convex and compact subset of E . A set valued map $A : D(A) \subseteq K \rightarrow E'$ is called monotone if whenever $v_i \in Au_i$, it follows that $\langle v_1 - v_2, u_1 - u_2 \rangle \geq 0$. The graph, denoted as $\mathcal{G}(A)$ consists of all pairs $[u, v]$ such that $v \in Au$. This is a monotone subset of $E \times E'$. Let $z \in E'$ be fixed. Show that

for $[u_i, v_i] \in \mathcal{G}(A)$, for $i = 1, 2, \dots, n$ there exists a solution $x \in K$ to the system of inequalities

$$\langle z + v_i, u_i - x \rangle \geq 0, i = 1, 2, \dots, n$$

Hint: Let P_n be all $\vec{\lambda} = (\lambda_1, \dots, \lambda_n)$ such that each $\lambda_k \geq 0$ and $\sum_{k=1}^n \lambda_k = 1$. Let $H : P_n \times P_n \rightarrow \mathbb{R}$ be given by

$$H(\vec{\mu}, \vec{\lambda}) \equiv \sum_{i=1}^n \mu_i \left\langle z + v_i, \sum_{j=1}^n \lambda_j u_j - u_i \right\rangle \quad (6.17)$$

Show that it is both convex and concave in both arguments. Then apply the min max theorem. Then argue that $H(\vec{\lambda}, \vec{\lambda}) \leq 0$ from monotonicity considerations. Letting $(\vec{\mu}_0, \vec{\lambda}_0)$ be the saddle point, you will have

$$\begin{aligned} H(\vec{\mu}, \vec{\lambda}_0) &\leq H(\vec{\mu}_0, \vec{\lambda}_0) \leq H(\vec{\mu}_0, \vec{\lambda}) \\ H(\vec{\mu}, \vec{\lambda}_0) &\leq H(\vec{\mu}_0, \vec{\lambda}_0) \leq H(\vec{\mu}_0, \vec{\mu}_0) \leq 0 \\ H(\vec{\mu}, \vec{\lambda}_0) &\leq 0 \end{aligned}$$

Now choose $\vec{\mu}$ judiciously while allowing $\vec{\lambda}_0$ to be used to define x which satisfies all the inequalities.

35. \uparrow It gets even better. Let $K_{u,v} \equiv \{x \in K : \langle z + v, u - x \rangle \geq 0\}$. Show that $K_{u,v}$ is compact and that the sets $K_{u,v}$ have the finite intersection property. Therefore, there exists $x \in \cap_{[u,v] \in \mathcal{G}(A)} K_{u,v}$. Explain why $\langle z + v, u - x \rangle \geq 0$ for all $[u, v] \in \mathcal{G}(A)$. What would the inequalities be if $-A$ were monotone?
36. Problem 33 gave a solution to the inequality $\langle f - z, y - x \rangle \leq 0, z \in Ax$ under the condition that A is upper semicontinuous. What are the differences between the result in the above problem and the result of Problem 33. You could replace A with $-A$ in the earlier problem. If you did, would you get the result of the above problem?
37. Are there convenient examples of monotone set valued maps? Yes, there are. Let X be a Banach space and let $\phi : X \rightarrow (-\infty, \infty]$ be convex, lower semicontinuous, and proper. See Problem 28 for a discussion of lower semicontinuous. Proper means that $\phi(x) < \infty$ for some x . Convex means the usual thing. $\phi(tx + (1-t)y) \leq t\phi(x) + (1-t)\phi(y)$ where $t \in [0, 1]$. Then $x^* \in \partial\phi(x)$ means that

$$\langle x^*, z - x \rangle \leq \phi(z) - \phi(x), \text{ for all } z \in X$$

Show that if $x^* \in \partial\phi(x)$, then $\phi(x) < \infty$. The set of points x where $\phi(x) < \infty$ is called the domain of ϕ denoted as $D(\phi)$. Also show that if $[x, x^*], [\hat{x}, \hat{x}^*]$ are two points of the graph of $\partial\phi$, then $\langle \hat{x}^* - x^*, \hat{x} - x \rangle \geq 0$ so that $\partial\phi$ is an example of a monotone graph. You might wonder whether this graph is nonempty. See the next problem for a partial answer to this question. Of course the above problem pertains to finite dimensional spaces so you could just take any $\phi : \mathbb{R}^n \rightarrow \mathbb{R}$ which is convex and differentiable. You can see that in this case the subgradient coincides with the derivative discussed later.

38. Let $\phi : \mathbb{R}^n \rightarrow \mathbb{R}$ be convex, proper lower semicontinuous, and bounded below. Show that the graph of $\partial\phi$ is nonempty. **Hint:** Just consider $\psi(x) = |x|^2 + \phi(x)$ and observe that this is coercive. Then argue using convexity that $\partial\psi(x) = \partial\phi(x) + 2x$. (You don't need to assume that ϕ is bounded below but it is convenient to assume this.)
39. Suppose $f : [0, T] \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ is continuous and an estimate of the following form holds. $\langle f(t, x), x \rangle \leq A + B|x|^2$ Show that there exists a solution to the initial value problem $x' = f(t, x)$, $x(0) = x_0$ for $t \in [0, T]$.
40. In the above problem, suppose that $-f + \alpha I$ is monotone for large enough α in addition to the estimate of that problem. Show that then there is only one solution to the problem. In fact, show that the solution depends continuously on the initial data.
41. It was shown that if $f : X \rightarrow X$ is locally Lipschitz where X is a Banach space. Then there exists a unique local solution to the IVP

$$y' = f(y), \quad y(0) = y_0$$

If f is bounded, then in fact the solutions exists on $[0, T]$ for any $T > 0$. Show that it suffices to assume that $\|f(y)\| \leq a + b\|y\|$.

42. Suppose $f(\cdot, \cdot) : \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ is continuous and also that $|f(t, x)| \leq M$ for all (t, x) . Show that there exists a solution to the initial value problem

$$x' = f(t, x), \quad x(0) = x_0 \in \mathbb{R}^n$$

for $t \in [0, T]$. **Hint:** You might consider $T : C([0, T], \mathbb{R}^n) \rightarrow C([0, T], \mathbb{R}^n)$ given by $Fx(t) \equiv x_0 + \int_0^t f(s, x(s)) ds$. Argue that F has a fixed point using the Schauder fixed point theorem.

43. Remove the assumption that $|f(t, x)| \leq M$ at the expense of obtaining only a local solution.
- Hint:** You can consider the closed set in \mathbb{R}^n $B = \overline{B(x_0, R)}$ where R is some positive number. Let P be the projection onto B .
44. In the Schauder fixed point theorem, eliminate the assumption that K is closed. **Hint:** You can argue that the $\{y_i\}$ in the approximation can be in $f(K)$.
45. Show that there is no one to one continuous function

$$f : [0, 1] \rightarrow \{(x, y) : x^2 + y^2 \leq 1\}$$

such that f is onto.

Chapter 7

The Derivative

7.1 Limits of a Function

As in the case of scalar valued functions of one variable, a concept closely related to continuity is that of the **limit of a function**. The notion of limit of a function makes sense at points x , which are limit points of $D(f)$ and this concept is defined next. In all that follows $(V, \|\cdot\|)$ and $(W, \|\cdot\|)$ are two normed linear spaces. Recall the definition of limit point first.

Definition 7.1.1 Let $A \subseteq W$ be a set. A point x , is a limit point of A if $B(x, r)$ contains infinitely many points of A for every $r > 0$.

Definition 7.1.2 Let $f : D(f) \subseteq V \rightarrow W$ be a function and let x be a limit point of $D(f)$. Then

$$\lim_{y \rightarrow x} f(y) = L$$

if and only if the following condition holds. For all $\varepsilon > 0$ there exists $\delta > 0$ such that if

$$0 < \|y - x\| < \delta, \text{ and } y \in D(f)$$

then,

$$\|L - f(y)\| < \varepsilon.$$

Theorem 7.1.3 If $\lim_{y \rightarrow x} f(y) = L$ and $\lim_{y \rightarrow x} f(y) = L_1$, then $L = L_1$.

Proof: Let $\varepsilon > 0$ be given. There exists $\delta > 0$ such that if $0 < \|y - x\| < \delta$ and $y \in D(f)$, then $\|f(y) - L\| < \varepsilon$, $\|f(y) - L_1\| < \varepsilon$. Pick such a y . There exists one because x is a limit point of $D(f)$. Then $\|L - L_1\| \leq \|L - f(y)\| + \|f(y) - L_1\| < \varepsilon + \varepsilon = 2\varepsilon$. Since $\varepsilon > 0$ was arbitrary, this shows $L = L_1$. ■

One can define what it means for $\lim_{y \rightarrow x} f(y) = \pm\infty$, as in the case of real valued functions.

Definition 7.1.4 If $f(x) \in \mathbb{R}$, $\lim_{y \rightarrow x} f(y) = \infty$ if for every number l , there exists $\delta > 0$ such that whenever $\|y - x\| < \delta$ and $y \in D(f)$, then $f(y) > l$. Also the assertion that $\lim_{y \rightarrow x} f(y) = -\infty$ means that for every number l , there exists $\delta > 0$ such that whenever $\|y - x\| < \delta$ and $y \in D(f)$, then $f(y) < l$.

The following theorem is just like the one variable version of calculus.

Theorem 7.1.5 Suppose $f : D(f) \subseteq V \rightarrow \mathbb{F}^m$. Then for x a limit point of $D(f)$,

$$\lim_{y \rightarrow x} f(y) = L \tag{7.1}$$

if and only if

$$\lim_{y \rightarrow x} f_k(y) = L_k \tag{7.2}$$

where $f(y) \equiv (f_1(y), \dots, f_p(y))$ and $L \equiv (L_1, \dots, L_p)$.

Suppose here that f has values in W , a normed linear space and

$$\lim_{y \rightarrow x} f(y) = L, \lim_{y \rightarrow x} g(y) = K$$

where $K, L \in W$. Then if $a, b \in \mathbb{F}$,

$$\lim_{y \rightarrow x} (af(y) + bg(y)) = aL + bK, \quad (7.3)$$

If W is an inner product space,

$$\lim_{y \rightarrow x} (f, g)(y) = (L, K) \quad (7.4)$$

If g is scalar valued with $\lim_{y \rightarrow x} g(y) = K$,

$$\lim_{y \rightarrow x} f(y)g(y) = LK. \quad (7.5)$$

Also, if h is a continuous function defined near L , then

$$\lim_{y \rightarrow x} h \circ f(y) = h(L). \quad (7.6)$$

Suppose $\lim_{y \rightarrow x} f(y) = L$. If $\|f(y) - b\| \leq r$ for all y sufficiently close to x , then $\|L - b\| \leq r$ also.

Proof: Suppose 7.1. Then letting $\varepsilon > 0$ be given there exists $\delta > 0$ such that if $0 < \|y - x\| < \delta$, it follows

$$\|f_k(y) - L_k\| \leq \|f(y) - L\| < \varepsilon$$

which verifies 7.2.

Now suppose 7.2 holds. Then letting $\varepsilon > 0$ be given, there exists δ_k such that if $0 < \|y - x\| < \delta_k$, then $\|f_k(y) - L_k\| < \varepsilon$. Let $0 < \delta < \min(\delta_1, \dots, \delta_p)$. Then if $0 < \|y - x\| < \delta$, it follows $\|f(y) - L\|_\infty < \varepsilon$. Any other norm on \mathbb{F}^m would work out the same way because the norms are all equivalent.

Each of the remaining assertions follows immediately from the coordinate descriptions of the various expressions and the first part. However, I will give a different argument for these.

The proof of 7.3 is left for you. Now 7.4 is to be verified. Let $\varepsilon > 0$ be given. Then by the triangle inequality,

$$\begin{aligned} |(f, g)(y) - (L, K)| &\leq |(f, g)(y) - (f(y), K)| + |(f(y), K) - (L, K)| \\ &\leq \|f(y)\| \|g(y) - K\| + \|K\| \|f(y) - L\|. \end{aligned}$$

There exists δ_1 such that if $0 < \|y - x\| < \delta_1$ and $y \in D(f)$, then $\|f(y) - L\| < 1$, and so for such y , the triangle inequality implies, $\|f(y)\| < 1 + \|L\|$. Therefore, for $0 < \|y - x\| < \delta_1$,

$$|(f, g)(y) - (L, K)| \leq (1 + \|K\| + \|L\|) [\|g(y) - K\| + \|f(y) - L\|]. \quad (7.7)$$

Now let $0 < \delta_2$ be such that if $y \in D(f)$ and $0 < \|x - y\| < \delta_2$,

$$\|f(y) - L\| < \frac{\varepsilon}{2(1 + \|K\| + \|L\|)}, \quad \|g(y) - K\| < \frac{\varepsilon}{2(1 + \|K\| + \|L\|)}.$$

Then letting $0 < \delta \leq \min(\delta_1, \delta_2)$, it follows from 7.7 that $|(f, g)(y) - (L, K)| < \varepsilon$ and this proves 7.4.

The proof of 7.5 is left to you.

Consider 7.6. Since h is continuous near L , it follows that for $\varepsilon > 0$ given, there exists $\eta > 0$ such that if $\|y - L\| < \eta$, then $\|h(y) - h(L)\| < \varepsilon$. Now since $\lim_{y \rightarrow x} f(y) = L$, there exists $\delta > 0$ such that if $0 < \|y - x\| < \delta$, then $\|f(y) - L\| < \eta$. Therefore, if $0 < \|y - x\| < \delta$, $\|h(f(y)) - h(L)\| < \varepsilon$.

It only remains to verify the last assertion. Assume $\|f(y) - b\| \leq r$. It is required to show that $\|L - b\| \leq r$. If this is not true, then $\|L - b\| > r$. Consider $B(L, \|L - b\| - r)$. Since L is the limit of f , it follows $f(y) \in B(L, \|L - b\| - r)$ whenever $y \in D(f)$ is close enough to x . Thus, by the triangle inequality, $\|f(y) - L\| < \|L - b\| - r$ and so

$$r < \|L - b\| - \|f(y) - L\| \leq \|b - L\| - \|f(y) - L\| \leq \|b - f(y)\|,$$

a contradiction to the assumption that $\|b - f(y)\| \leq r$. ■

The relation between continuity and limits is as follows.

Theorem 7.1.6 For $f : D(f) \rightarrow W$ and $x \in D(f)$ a limit point of $D(f)$, f is continuous at x if and only if $\lim_{y \rightarrow x} f(y) = f(x)$.

Proof: First suppose f is continuous at x a limit point of $D(f)$. Then for every $\varepsilon > 0$ there exists $\delta > 0$ such that if $\|x - y\| < \delta$ and $y \in D(f)$, then $|f(x) - f(y)| < \varepsilon$. In particular, this holds if $0 < \|x - y\| < \delta$ and this is just the definition of the limit. Hence $f(x) = \lim_{y \rightarrow x} f(y)$.

Next suppose x is a limit point of $D(f)$ and $\lim_{y \rightarrow x} f(y) = f(x)$. This means that if $\varepsilon > 0$ there exists $\delta > 0$ such that for $0 < \|x - y\| < \delta$ and $y \in D(f)$, it follows $|f(y) - f(x)| < \varepsilon$. However, if $y = x$, then $|f(y) - f(x)| = |f(x) - f(x)| = 0$ and so whenever $y \in D(f)$ and $\|x - y\| < \delta$, it follows $|f(x) - f(y)| < \varepsilon$, showing f is continuous at x . ■

Example 7.1.7 Find $\lim_{(x,y) \rightarrow (3,1)} \left(\frac{x^2-9}{x-3}, y \right)$.

It is clear that $\lim_{(x,y) \rightarrow (3,1)} \frac{x^2-9}{x-3} = 6$ and $\lim_{(x,y) \rightarrow (3,1)} y = 1$. Therefore, this limit equals $(6, 1)$.

Example 7.1.8 Find $\lim_{(x,y) \rightarrow (0,0)} \frac{xy}{x^2+y^2}$.

First of all, observe the domain of the function is $\mathbb{R}^2 \setminus \{(0,0)\}$, every point in \mathbb{R}^2 except the origin. Therefore, $(0,0)$ is a limit point of the domain of the function so it might make sense to take a limit. However, just as in the case of a function of one variable, the limit may not exist. In fact, this is the case here. To see this, take points on the line $y = 0$. At these points, the value of the function equals 0. Now consider points on the line $y = x$ where the value of the function equals $1/2$. Since, arbitrarily close to $(0,0)$, there are points where the function equals $1/2$ and points where the function has the value 0, it follows there can be no limit. Just take $\varepsilon = 1/10$ for example. You cannot be within $1/10$ of $1/2$ and also within $1/10$ of 0 at the same time.

Note it is necessary to rely on the definition of the limit much more than in the case of a function of one variable and there are no easy ways to do limit problems for functions of more than one variable. It is what it is and you will not deal with these concepts without suffering and anguish.

7.2 Basic Definitions

The concept of derivative generalizes right away to functions of many variables. However, no attempt will be made to consider derivatives from one side or another. This is because when you consider functions of many variables, there isn't a well defined side. However, it is certainly the case that there are more general notions which include such things. I will present a fairly general notion of the derivative of a function which is defined on a normed vector space which has values in a normed vector space. The case of most interest is that of a function which maps \mathbb{F}^n to \mathbb{F}^m but it is no more trouble to consider the extra generality and it is sometimes useful to have this extra generality because sometimes you want to consider functions defined, for example on subspaces of \mathbb{F}^n and it is nice to not have to trouble with ad hoc considerations. Also, you might want to consider \mathbb{F}^n with some norm other than the usual one.

In what follows, X, Y will denote normed vector spaces. Thanks to Theorem 5.2.4 all the definitions and theorems given below work the same for any norm given on the vector spaces.

Let U be an open set in X , and let $f : U \rightarrow Y$ be a function.

Definition 7.2.1 A function g is $o(v)$ if

$$\lim_{\|v\| \rightarrow 0} \frac{g(v)}{\|v\|} = 0 \quad (7.8)$$

A function $f : U \rightarrow Y$ is differentiable at $x \in U$ if there exists a linear transformation $L \in \mathcal{L}(X, Y)$ such that

$$f(x + v) = f(x) + Lv + o(v)$$

This linear transformation L is the definition of $Df(x)$. This derivative is often called the Frechet derivative.

Note that from Theorem 5.2.4 the question whether a given function is differentiable is independent of the norm used on the finite dimensional vector space. That is, a function is differentiable with one norm if and only if it is differentiable with another norm.

The definition 7.8 means the error $f(x + v) - f(x) - Lv$ converges to 0 faster than $\|v\|$. Thus the above definition is equivalent to saying

$$\lim_{\|v\| \rightarrow 0} \frac{\|f(x + v) - f(x) - Lv\|}{\|v\|} = 0 \quad (7.9)$$

or equivalently,

$$\lim_{y \rightarrow x} \frac{\|f(y) - f(x) - Df(x)(y - x)\|}{\|y - x\|} = 0. \quad (7.10)$$

The symbol, $o(v)$ should be thought of as an adjective. Thus, if t and k are constants,

$$o(v) = o(v) + o(v), \quad o(tv) = o(v), \quad ko(v) = o(v)$$

and other similar observations hold.

Theorem 7.2.2 The derivative is well defined.

Proof: First note that for a fixed nonzero vector \mathbf{v} , $\mathbf{o}(t\mathbf{v}) = \mathbf{o}(t)$. This is because

$$\lim_{t \rightarrow 0} \frac{\mathbf{o}(t\mathbf{v})}{|t|} = \lim_{t \rightarrow 0} \|\mathbf{v}\| \frac{\mathbf{o}(t\mathbf{v})}{\|t\mathbf{v}\|} = \mathbf{0}$$

Now suppose both L_1 and L_2 work in the above definition. Then let \mathbf{v} be any vector and let t be a real scalar which is chosen small enough that $t\mathbf{v} + \mathbf{x} \in U$. Then

$$\mathbf{f}(\mathbf{x} + t\mathbf{v}) = \mathbf{f}(\mathbf{x}) + L_1 t\mathbf{v} + \mathbf{o}(t\mathbf{v}), \quad \mathbf{f}(\mathbf{x} + t\mathbf{v}) = \mathbf{f}(\mathbf{x}) + L_2 t\mathbf{v} + \mathbf{o}(t\mathbf{v}).$$

Therefore, subtracting these two yields $(L_2 - L_1)(t\mathbf{v}) = \mathbf{o}(t\mathbf{v}) = \mathbf{o}(t)$. Therefore, dividing by t yields $(L_2 - L_1)(\mathbf{v}) = \frac{\mathbf{o}(t)}{t}$. Now let $t \rightarrow 0$ to conclude that $(L_2 - L_1)(\mathbf{v}) = \mathbf{0}$. Since this is true for all \mathbf{v} , it follows $L_2 = L_1$. This proves the theorem. ■

In the following lemma, $\|D\mathbf{f}(\mathbf{x})\|$ is the operator norm of the linear transformation, $D\mathbf{f}(\mathbf{x})$.

Lemma 7.2.3 *Let \mathbf{f} be differentiable at \mathbf{x} . Then \mathbf{f} is continuous at \mathbf{x} and in fact, there exists $K > 0$ such that whenever $\|\mathbf{v}\|$ is small enough,*

$$\|\mathbf{f}(\mathbf{x} + \mathbf{v}) - \mathbf{f}(\mathbf{x})\| \leq K \|\mathbf{v}\|$$

Also if \mathbf{f} is differentiable at \mathbf{x} , then

$$\mathbf{o}(\|\mathbf{f}(\mathbf{x} + \mathbf{v}) - \mathbf{f}(\mathbf{x})\|) = \mathbf{o}(\mathbf{v})$$

Proof: From the definition of the derivative,

$$\mathbf{f}(\mathbf{x} + \mathbf{v}) - \mathbf{f}(\mathbf{x}) = D\mathbf{f}(\mathbf{x})\mathbf{v} + \mathbf{o}(\mathbf{v}).$$

Let $\|\mathbf{v}\|$ be small enough that $\frac{\mathbf{o}(\|\mathbf{v}\|)}{\|\mathbf{v}\|} < 1$ so that $\|\mathbf{o}(\mathbf{v})\| \leq \|\mathbf{v}\|$. Then for such \mathbf{v} ,

$$\|\mathbf{f}(\mathbf{x} + \mathbf{v}) - \mathbf{f}(\mathbf{x})\| \leq \|D\mathbf{f}(\mathbf{x})\mathbf{v}\| + \|\mathbf{v}\| \leq (\|D\mathbf{f}(\mathbf{x})\| + 1)\|\mathbf{v}\|$$

This proves the lemma with $K = \|D\mathbf{f}(\mathbf{x})\| + 1$. Recall the operator norm discussed in Definition 5.2.2.

The last assertion is implied by the first as follows. Define

$$\mathbf{h}(\mathbf{v}) \equiv \begin{cases} \frac{\mathbf{o}(\|\mathbf{f}(\mathbf{x} + \mathbf{v}) - \mathbf{f}(\mathbf{x})\|)}{\|\mathbf{f}(\mathbf{x} + \mathbf{v}) - \mathbf{f}(\mathbf{x})\|} & \text{if } \|\mathbf{f}(\mathbf{x} + \mathbf{v}) - \mathbf{f}(\mathbf{x})\| \neq 0 \\ \mathbf{0} & \text{if } \|\mathbf{f}(\mathbf{x} + \mathbf{v}) - \mathbf{f}(\mathbf{x})\| = 0 \end{cases}$$

Then $\lim_{\|\mathbf{v}\| \rightarrow 0} \mathbf{h}(\mathbf{v}) = \mathbf{0}$ from continuity of \mathbf{f} at \mathbf{x} which is implied by the first part. Also from the above estimate, if $\|\mathbf{v}\|$ is sufficiently small,

$$\left\| \frac{\mathbf{o}(\|\mathbf{f}(\mathbf{x} + \mathbf{v}) - \mathbf{f}(\mathbf{x})\|)}{\|\mathbf{v}\|} \right\| = \|\mathbf{h}(\mathbf{v})\| \frac{\|\mathbf{f}(\mathbf{x} + \mathbf{v}) - \mathbf{f}(\mathbf{x})\|}{\|\mathbf{v}\|} \leq \|\mathbf{h}(\mathbf{v})\| (\|D\mathbf{f}(\mathbf{x})\| + 1)$$

and $\lim_{\|\mathbf{v}\| \rightarrow 0} \|\mathbf{h}(\mathbf{v})\| = 0$. This establishes the second claim. ■

7.3 The Chain Rule

With the above lemma, it is easy to prove the chain rule.

Theorem 7.3.1 (The chain rule) *Let U and V be open sets $U \subseteq X$ and $V \subseteq Y$. Suppose $f : U \rightarrow V$ is differentiable at $x \in U$ and suppose $g : V \rightarrow \mathbb{R}^q$ is differentiable at $f(x) \in V$. Then $g \circ f$ is differentiable at x and*

$$D(g \circ f)(x) = Dg(f(x))Df(x).$$

Proof: This follows from a computation. Let $B(x, r) \subseteq U$ and let r also be small enough that for $\|v\| \leq r$, it follows that $f(x + v) \in V$. Such an r exists because f is continuous at x . For $\|v\| < r$, the definition of differentiability of g and f implies

$$\begin{aligned} g(f(x + v)) - g(f(x)) &= \\ &= Dg(f(x))(f(x + v) - f(x)) + o(f(x + v) - f(x)) \\ &= Dg(f(x))[Df(x)v + o(v)] + o(f(x + v) - f(x)) \\ &= D(g(f(x)))D(f(x))v + o(v) + o(f(x + v) - f(x)) \\ &= D(g(f(x)))D(f(x))v + o(v) \end{aligned} \quad (7.11)$$

By Lemma 7.2.3. From the definition of the derivative $D(g \circ f)(x)$ exists and equals $D(g(f(x)))D(f(x))$. ■

7.4 The Matrix of the Derivative

The case of most interest here is the only one I will discuss. It is the case where $X = \mathbb{R}^n$ and $Y = \mathbb{R}^m$, the function being defined on an open subset of \mathbb{R}^n . Of course this all generalizes to arbitrary vector spaces and one considers the matrix taken with respect to various bases. However, I am going to restrict to the case just mentioned here. As above, f will be defined and differentiable on an open set $U \subseteq \mathbb{R}^n$.

As discussed in the review material on linear maps, the matrix of $Df(x)$ is the matrix having the i^{th} column equal to $Df(x)e_i$ and so it is only necessary to compute this. Let t be a small real number such that

$$\frac{f(x + te_i) - f(x) - Df(x)(te_i)}{t} = \frac{o(t)}{t}$$

Therefore,

$$\frac{f(x + te_i) - f(x)}{t} = Df(x)(e_i) + \frac{o(t)}{t}$$

The limit exists on the right and so it exists on the left also. Thus

$$\frac{\partial f(x)}{\partial x_i} \equiv \lim_{t \rightarrow 0} \frac{f(x + te_i) - f(x)}{t} = Df(x)(e_i)$$

and so the matrix of the derivative is just the matrix which has the i^{th} column equal to the i^{th} partial derivative of f . Note that this shows that whenever f is differentiable, it follows that the partial derivatives all exist. It does not go the other way however as discussed later.

Theorem 7.4.1 Let $\mathbf{f} : U \subseteq \mathbb{F}^n \rightarrow \mathbb{F}^m$ and suppose \mathbf{f} is differentiable at \mathbf{x} . Then all the partial derivatives $\frac{\partial f_i(\mathbf{x})}{\partial x_j}$ exist and if $\mathbf{J}\mathbf{f}(\mathbf{x})$ is the matrix of the linear transformation, $D\mathbf{f}(\mathbf{x})$ with respect to the standard basis vectors, then the ij^{th} entry is given by $\frac{\partial f_i}{\partial x_j}(\mathbf{x})$ also denoted as $f_{i,j}$ or f_{i,x_j} . It is the matrix whose i^{th} column is

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_i} \equiv \lim_{t \rightarrow 0} \frac{\mathbf{f}(\mathbf{x} + t\mathbf{e}_i) - \mathbf{f}(\mathbf{x})}{t}.$$

Of course there is a generalization of this idea called the directional derivative.

Definition 7.4.2 In general, the symbol $D_v \mathbf{f}(\mathbf{x})$ is defined by

$$\lim_{t \rightarrow 0} \frac{\mathbf{f}(\mathbf{x} + t\mathbf{v}) - \mathbf{f}(\mathbf{x})}{t}$$

where $t \in \mathbb{F}$. In case $|\mathbf{v}| = 1$, $\mathbb{F} = \mathbb{R}$, and the norm is the standard Euclidean norm, this is called the directional derivative. More generally, with no restriction on the size of \mathbf{v} and in any linear space, it is called the Gateaux derivative. \mathbf{f} is said to be Gateaux differentiable at \mathbf{x} if there exists $D_v \mathbf{f}(\mathbf{x})$ such that

$$\lim_{t \rightarrow 0} \frac{\mathbf{f}(\mathbf{x} + t\mathbf{v}) - \mathbf{f}(\mathbf{x})}{t} = D_v \mathbf{f}(\mathbf{x})$$

where $\mathbf{v} \rightarrow D_v \mathbf{f}(\mathbf{x})$ is linear. Thus we say it is Gateaux differentiable if the Gateaux derivative exists for each \mathbf{v} and $\mathbf{v} \rightarrow D_v \mathbf{f}(\mathbf{x})$ is linear. Note that $\frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_i} = D_{\mathbf{e}_i} \mathbf{f}(\mathbf{x})$.¹

What if all the partial derivatives of \mathbf{f} exist? Does it follow that \mathbf{f} is differentiable? Consider the following function, $f : \mathbb{R}^2 \rightarrow \mathbb{R}$,

$$f(x, y) = \begin{cases} \frac{xy}{x^2 + y^2} & \text{if } (x, y) \neq (0, 0) \\ 0 & \text{if } (x, y) = (0, 0) \end{cases}.$$

Then from the definition of partial derivatives,

$$\lim_{h \rightarrow 0} \frac{f(h, 0) - f(0, 0)}{h} = \lim_{h \rightarrow 0} \frac{0 - 0}{h} = 0$$

and

$$\lim_{h \rightarrow 0} \frac{f(0, h) - f(0, 0)}{h} = \lim_{h \rightarrow 0} \frac{0 - 0}{h} = 0$$

However f is not even continuous at $(0, 0)$ which may be seen by considering the behavior of the function along the line $y = x$ and along the line $x = 0$. By Lemma 7.2.3 this implies f is not differentiable. Therefore, it is necessary to consider the correct definition of the derivative given above if you want to get a notion which generalizes the concept of the derivative of a function of one variable in such a way as to preserve continuity whenever the function is differentiable.

What if the one dimensional derivative in the definition of the Gateaux derivative exists for all nonzero \mathbf{v} ? Is the function differentiable then? Maybe not. See Problem 12 in the exercises for example.

¹René Gateaux was one of the many young French men killed in world war I. This derivative is named after him, but it developed naturally from ideas used in the calculus of variations which were due to Euler and Lagrange back in the 1700's.

7.5 A Mean Value Inequality

The following theorem will be very useful in much of what follows. It is a version of the mean value theorem as is the next lemma. The mean value theorem depends on the function having values in \mathbb{R} and in the lemma and theorem, it has values in a normed vector space.

Lemma 7.5.1 *Let Y be a normed vector space and suppose $\mathbf{h} : [0, 1] \rightarrow Y$ is continuous and differentiable from the right and satisfies $\|\mathbf{h}'(t)\| \leq M$, $M \geq 0$. Then $\|\mathbf{h}(1) - \mathbf{h}(0)\| \leq M$.*

Proof: Let $\varepsilon > 0$ be given and let

$$S \equiv \{t \in [0, 1] : \text{for all } s \in [0, t], \|\mathbf{h}(s) - \mathbf{h}(0)\| \leq (M + \varepsilon)s\}$$

Then $0 \in S$. Let $t = \sup S$. Then by continuity of \mathbf{h} it follows

$$\|\mathbf{h}(t) - \mathbf{h}(0)\| = (M + \varepsilon)t \quad (7.12)$$

Suppose $t < 1$. Then there exist positive numbers, h_k decreasing to 0 such that

$$\|\mathbf{h}(t + h_k) - \mathbf{h}(0)\| > (M + \varepsilon)(t + h_k)$$

and now it follows from 7.12 and the triangle inequality that

$$\begin{aligned} & \|\mathbf{h}(t + h_k) - \mathbf{h}(t)\| + \|\mathbf{h}(t) - \mathbf{h}(0)\| \\ = & \|\mathbf{h}(t + h_k) - \mathbf{h}(t)\| + (M + \varepsilon)t > (M + \varepsilon)(t + h_k) \end{aligned}$$

Thus

$$\|\mathbf{h}(t + h_k) - \mathbf{h}(t)\| > (M + \varepsilon)h_k$$

Now dividing by h_k and letting $k \rightarrow \infty$, $\|\mathbf{h}'(t)\| \geq M + \varepsilon$, a contradiction. Thus $t = 1$. Since ε is arbitrary, the conclusion of the lemma follows. ■

Theorem 7.5.2 *Suppose U is an open subset of X and $\mathbf{f} : U \rightarrow Y$ has the property that $D\mathbf{f}(\mathbf{x})$ exists for all \mathbf{x} in U and that, $\mathbf{x} + t(\mathbf{y} - \mathbf{x}) \in U$ for all $t \in [0, 1]$. (The line segment joining the two points lies in U .) Suppose also that for all points on this line segment, $\|D\mathbf{f}(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))\| \leq M$. Then $\|\mathbf{f}(\mathbf{y}) - \mathbf{f}(\mathbf{x})\| \leq M\|\mathbf{y} - \mathbf{x}\|$. More generally if $\|D_{\mathbf{v}}\mathbf{f}(\mathbf{y})\| \leq M$ for all \mathbf{y} on the segment joining \mathbf{x} and $\mathbf{x} + \mathbf{v}$, then $\|\mathbf{f}(\mathbf{x} + a\mathbf{v}) - \mathbf{f}(\mathbf{x})\| \leq Ma$. Also $D_{a\mathbf{v}}\mathbf{f}(\mathbf{x}) = aD_{\mathbf{v}}\mathbf{f}(\mathbf{x})$ if $a \neq 0$.*

Proof: Let $\mathbf{h}(t) \equiv \mathbf{f}(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))$. Then by the chain rule applied to $\mathbf{h}(t)$, $\mathbf{h}'(t) = D\mathbf{f}(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))(\mathbf{y} - \mathbf{x})$ and so

$$\|\mathbf{h}'(t)\| = \|D\mathbf{f}(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))(\mathbf{y} - \mathbf{x})\| \leq M\|\mathbf{y} - \mathbf{x}\|$$

by Lemma 7.5.1, $\|\mathbf{h}(1) - \mathbf{h}(0)\| = \|\mathbf{f}(\mathbf{y}) - \mathbf{f}(\mathbf{x})\| \leq M\|\mathbf{y} - \mathbf{x}\|$. For the second part, let $\mathbf{h}(t) \equiv \mathbf{f}(\mathbf{x} + t\mathbf{v})$. Then if $a \neq 0$,

$$\begin{aligned} \mathbf{h}'(t) &= \lim_{h \rightarrow 0} \frac{\mathbf{h}(t+h) - \mathbf{h}(t)}{h} \equiv \lim_{h \rightarrow 0} \frac{a}{ha} (\mathbf{f}(\mathbf{x} + t\mathbf{v} + h\mathbf{v}) - \mathbf{f}(\mathbf{x} + t\mathbf{v})) \\ &= D_{\mathbf{v}}\mathbf{f}(\mathbf{x} + t\mathbf{v})a. \end{aligned}$$

This shows that $D_{a\mathbf{v}}\mathbf{f}(\mathbf{x}) = aD_{\mathbf{v}}\mathbf{f}(\mathbf{x})$. Now for the inequality, there is nothing to show if $a = 0$ so assume $a \neq 0$. Then by assumption and Lemma 7.5.1, $\|\mathbf{h}(1) - \mathbf{h}(0)\| = \|\mathbf{f}(\mathbf{x} + a\mathbf{v}) - \mathbf{f}(\mathbf{x})\| \leq Ma$. ■

7.6 Existence of the Derivative, C^1 Functions

There is a way to get the differentiability of a function from the existence and continuity of one dimensional directional derivatives. The following theorem is the main result. It gives easy to verify one dimensional conditions for the existence of the derivative. The meaning of $\|\cdot\|$ will be determined by context in what follows. This theorem says that if the Gateaux derivatives exist for each vector in a basis and they are also continuous, then the function is differentiable.

Theorem 7.6.1 *Let X be a normed vector space having basis $\{v_1, \dots, v_n\}$ and let Y be another normed vector space. Let U be an open set in X and let $f : U \rightarrow Y$ have the property that the one dimensional limits*

$$D_{v_k} f(x) \equiv \lim_{t \rightarrow 0} \frac{f(x + tv_k) - f(x)}{t}$$

exist and $x \rightarrow D_{v_k} f(x)$ are continuous functions of $x \in U$ as functions with values in Y . Then $Df(x)$ exists and

$$Df(x) v = \sum_{k=1}^n D_{v_k} f(x) a_k$$

where $v = \sum_{k=1}^n a_k v_k$. Furthermore, $x \rightarrow Df(x)$ is continuous; that is

$$\lim_{y \rightarrow x} \|Df(y) - Df(x)\| = 0.$$

Proof: Let $v = \sum_{k=1}^n a_k v_k$ where all a_k are small enough that for all $k \geq 0$,

$$x + \sum_{j=1}^k a_j v_j \in \overline{B(x, r)} \subseteq U, \sum_{k=1}^0 a_k v_k \equiv 0.$$

The mapping $v \rightarrow (a_1, \dots, a_n)$ is an isomorphism of V and \mathbb{F}^n and we can define a norm as $\sum_k |a_k|$ which is equivalent to the norm on V thanks to Theorem 5.2.4. Let $h_k(x) \equiv f(x + \sum_{j=1}^{k-1} a_j v_j) - f(x)$. Then collecting the terms,

$$f(x + v) - f(x) = \sum_{k=1}^n (h_k(x + a_k v_k) - h_k(x)) + \sum_{k=1}^n (f(x + a_k v_k) - f(x)) \quad (7.13)$$

Using Theorem 7.5.2,

$$\begin{aligned} \|D_{a_k v_k} h_k(x + ta_k v_k)\| &= \|a_k D_{v_k} h_k(x + ta_k v_k)\| \\ &= \left\| a_k \left(D_{v_k} f \left(x + \sum_{j=1}^{k-1} a_j v_j + ta_k v_k \right) - D_{v_k} f(x + ta_k v_k) \right) \right\| \\ &\leq C \|v\| \varepsilon \end{aligned}$$

provided $\|v\|$ is sufficiently small, thanks to the assumption that the $D_{v_k} f$ are continuous. It follows, since ε is arbitrary that the first sum on the right in 7.13 is $o(v)$. Now

$$(f(x + a_k v_k) - f(x)) - D_{v_k} f(x) a_k =$$

$$\mathbf{f}(\mathbf{x} + a_k \mathbf{v}_k) - (\mathbf{f}(\mathbf{x}) + D_{\mathbf{v}_k} \mathbf{f}(\mathbf{x}) a_k) = a_k \left(\frac{\mathbf{f}(\mathbf{x} + a_k \mathbf{v}_k) - \mathbf{f}(\mathbf{x})}{a_k} - D_{\mathbf{v}_k} \mathbf{f}(\mathbf{x}) \right) = \mathbf{o}(v)$$

because

$$\begin{aligned} & \left\| a_k \left(\frac{\mathbf{f}(\mathbf{x} + a_k \mathbf{v}_k) - \mathbf{f}(\mathbf{x})}{a_k} - D_{\mathbf{v}_k} \mathbf{f}(\mathbf{x}) \right) \right\| \\ & \leq \| \mathbf{v} \| \left\| \left(\frac{\mathbf{f}(\mathbf{x} + a_k \mathbf{v}_k) - \mathbf{f}(\mathbf{x})}{a_k} - D_{\mathbf{v}_k} \mathbf{f}(\mathbf{x}) \right) \right\|. \end{aligned}$$

Collecting terms in 7.13,

$$\mathbf{f}(\mathbf{x} + \mathbf{v}) - \mathbf{f}(\mathbf{x}) = \mathbf{o}(v) + \sum_{k=1}^n (\mathbf{f}(\mathbf{x} + a_k \mathbf{v}_k) - \mathbf{f}(\mathbf{x})) = \mathbf{o}(v) + \sum_{k=1}^n D_{\mathbf{v}_k} \mathbf{f}(\mathbf{x}) a_k$$

which shows that $D\mathbf{f}(\mathbf{x})(\mathbf{v}) = \sum_{k=1}^n D_{\mathbf{v}_k} \mathbf{f}(\mathbf{x}) a_k$ where $\mathbf{v} = \sum_{k=1}^n a_k \mathbf{v}_k$. This formula also shows that $\mathbf{x} \rightarrow D\mathbf{f}(\mathbf{x})$ is continuous because of the continuity of these $D_{\mathbf{v}_k} \mathbf{f}$. ■

Note how if $X = \mathbb{R}^p$ and the basis vectors are the \mathbf{e}_k , then the \mathbf{a} are just the components of the vector \mathbf{v} taken with respect to the usual basis vectors. Thus this gives the above result about the matrix of $D\mathbf{f}(\mathbf{x})$.

This motivates the following definition of what it means for a function to be C^1 .

Definition 7.6.2 Let U be an open subset of a normed finite dimensional vector space, X and let $\mathbf{f} : U \rightarrow Y$ another finite dimensional normed vector space. Then \mathbf{f} is said to be C^1 if there exists a basis for $X, \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ such that the Gateaux derivatives, $D_{\mathbf{v}_k} \mathbf{f}(\mathbf{x})$ exist on U and are continuous functions of \mathbf{x} .

Note that as a special case where $X = \mathbb{R}^n$, you could let the $\mathbf{v}_k = \mathbf{e}_k$ and the condition would reduce to nothing more than a statement that the partial derivatives $\frac{\partial \mathbf{f}}{\partial x_i}$ are all continuous. If $X = \mathbb{R}$, this is not a very interesting condition. It would say the derivative exists if the derivative exists and is continuous.

Here is another definition of what it means for a function to be C^1 .

Definition 7.6.3 Let U be an open subset of a normed finite dimensional vector space, X and let $\mathbf{f} : U \rightarrow Y$ another finite dimensional normed vector space. Then \mathbf{f} is said to be C^1 if \mathbf{f} is differentiable and $\mathbf{x} \rightarrow D\mathbf{f}(\mathbf{x})$ is continuous as a map from U to $\mathcal{L}(X, Y)$.

Now the following major theorem states these two definitions are equivalent. This is obviously so in the special case where $X = \mathbb{R}^n$ and the special basis is the usual one because, as observed earlier, the matrix of $D\mathbf{f}(\mathbf{x})$ is just the one which has for its columns the partial derivatives which are given to be continuous.

Theorem 7.6.4 Let U be an open subset of a normed finite dimensional vector space X and let $\mathbf{f} : U \rightarrow Y$ another finite dimensional normed vector space. Then the two definitions above are equivalent.

Proof: It was shown in Theorem 7.6.1, the one about the continuity of the Gateaux derivatives yielding differentiability, that Definition 7.6.2 implies 7.6.3. Suppose then that Definition 7.6.3 holds. Then if \mathbf{v} is any vector,

$$\lim_{t \rightarrow 0} \frac{\mathbf{f}(\mathbf{x} + t\mathbf{v}) - \mathbf{f}(\mathbf{x})}{t} = \lim_{t \rightarrow 0} \frac{D\mathbf{f}(\mathbf{x})t\mathbf{v} + \mathbf{o}(t\mathbf{v})}{t} = D\mathbf{f}(\mathbf{x})\mathbf{v} + \lim_{t \rightarrow 0} \frac{\mathbf{o}(t\mathbf{v})}{t} = D\mathbf{f}(\mathbf{x})\mathbf{v}$$

Thus $D_v \mathbf{f}(x)$ exists and equals $D\mathbf{f}(x)v$. By continuity of $x \rightarrow D\mathbf{f}(x)$, this establishes continuity of $x \rightarrow D_v \mathbf{f}(x)$ and proves the theorem. ■

Note that the proof of the theorem also implies the following corollary.

Corollary 7.6.5 *Let U be an open subset of a normed finite dimensional vector space, X and let $\mathbf{f} : U \rightarrow Y$ another finite dimensional normed vector space. Then if there is a basis of X , $\{v_1, \dots, v_n\}$ such that the Gateaux derivatives, $D_{v_k} \mathbf{f}(x)$ exist and are continuous, then all Gateaux derivatives, $D_v \mathbf{f}(x)$ exist and are continuous for all $v \in X$. Also $D\mathbf{f}(x)(v) = D_v \mathbf{f}(x)$.*

From now on, whichever definition is more convenient will be used.

7.7 Higher Order Derivatives

If $\mathbf{f} : U \subseteq X \rightarrow Y$ for U an open set, then $x \rightarrow D\mathbf{f}(x)$ is a mapping from U to $\mathcal{L}(X, Y)$, a normed vector space. Therefore, it makes perfect sense to ask whether this function is also differentiable.

Definition 7.7.1 *The following is the definition of the second derivative. $D^2 \mathbf{f}(x) \equiv D(D\mathbf{f}(x))$.*

Thus, $D\mathbf{f}(x+v) - D\mathbf{f}(x) = D^2 \mathbf{f}(x)v + o(v)$. This implies

$$D^2 \mathbf{f}(x) \in \mathcal{L}(X, \mathcal{L}(X, Y)), \quad D^2 \mathbf{f}(x)(u)(v) \in Y,$$

and the map $(u, v) \rightarrow D^2 \mathbf{f}(x)(u)(v)$ is a bilinear map having values in Y . In other words, the two functions,

$$u \rightarrow D^2 \mathbf{f}(x)(u)(v), \quad v \rightarrow D^2 \mathbf{f}(x)(u)(v)$$

are both linear.

The same pattern applies to taking higher order derivatives. For example, $D^3 \mathbf{f}(x) \equiv D(D^2 \mathbf{f}(x))$ and $D^3 \mathbf{f}(x)$ may be considered as a trilinear map having values in Y . In general $D^k \mathbf{f}(x)$ may be considered a k linear map. This means

$$(u_1, \dots, u_k) \rightarrow D^k \mathbf{f}(x)(u_1) \cdots (u_k)$$

has the property $u_j \rightarrow D^k \mathbf{f}(x)(u_1) \cdots (u_j) \cdots (u_k)$ is linear.

Also, instead of writing $D^2 \mathbf{f}(x)(u)(v)$, or $D^3 \mathbf{f}(x)(u)(v)(w)$ the following notation is often used.

$$D^2 \mathbf{f}(x)(u, v) \text{ or } D^3 \mathbf{f}(x)(u, v, w)$$

with similar conventions for higher derivatives than 3. Another convention which is often used is the notation $D^k \mathbf{f}(x)v^k$ instead of $D^k \mathbf{f}(x)(v, \dots, v)$.

Note that for every k , $D^k \mathbf{f}$ maps U to a normed vector space. As mentioned above, $D\mathbf{f}(x)$ has values in $\mathcal{L}(X, Y)$, $D^2 \mathbf{f}(x)$ has values in $\mathcal{L}(X, \mathcal{L}(X, Y))$, etc. Thus it makes sense to consider whether $D^k \mathbf{f}$ is continuous. This is described in the following definition.

Definition 7.7.2 *Let U be an open subset of X , a normed vector space, and let $\mathbf{f} : U \rightarrow Y$. Then \mathbf{f} is $C^k(U)$ if \mathbf{f} and its first k derivatives are all continuous. Also, $D^k \mathbf{f}(x)$ when it exists can be considered a Y valued multi-linear function. Sometimes these are called tensors in case \mathbf{f} has scalar values.*

7.8 Some Standard Notation

In the case where $X = \mathbb{R}^n$ there is a special notation which is often used to describe higher order mixed partial derivatives. It is called multi-index notation.

Definition 7.8.1 $\alpha = (\alpha_1, \dots, \alpha_n)$ for $\alpha_1 \cdots \alpha_n$ positive integers is called a multi-index, as before with polynomials. For α a multi-index, $|\alpha| \equiv \alpha_1 + \cdots + \alpha_n$, and if $\mathbf{x} \in X$,

$$\mathbf{x} = (x_1, \dots, x_n),$$

and \mathbf{f} a function, define

$$\mathbf{x}^\alpha \equiv x_1^{\alpha_1} x_2^{\alpha_2} \cdots x_n^{\alpha_n}, \quad D^\alpha \mathbf{f}(\mathbf{x}) \equiv \frac{\partial^{|\alpha|} \mathbf{f}(\mathbf{x})}{\partial x_1^{\alpha_1} \partial x_2^{\alpha_2} \cdots \partial x_n^{\alpha_n}}.$$

Then in this special case, the following is another description of what is meant by a C^k function.

Definition 7.8.2 Let U be an open subset of \mathbb{R}^n and let $\mathbf{f} : U \rightarrow Y$. Then for k a nonnegative integer, a differentiable function \mathbf{f} is C^k if for every $|\alpha| \leq k$, $D^\alpha \mathbf{f}$ exists and is continuous.

Theorem 7.8.3 Let U be an open subset of \mathbb{R}^n and let $\mathbf{f} : U \rightarrow Y$. Then if $D^r \mathbf{f}(\mathbf{x})$ exists for $r \leq k$, then $D^r \mathbf{f}$ is continuous at \mathbf{x} for $r \leq k$ if and only if $D^\alpha \mathbf{f}$ is continuous at \mathbf{x} for each $|\alpha| \leq k$.

Proof: First consider the case of a single derivative. Then as shown above, the matrix of $D\mathbf{f}(\mathbf{x})$ is just

$$J(\mathbf{x}) \equiv \begin{pmatrix} \frac{\partial \mathbf{f}}{\partial x_1}(\mathbf{x}) & \cdots & \frac{\partial \mathbf{f}}{\partial x_n}(\mathbf{x}) \end{pmatrix}$$

and to say that $\mathbf{x} \rightarrow D\mathbf{f}(\mathbf{x})$ is continuous is the same as saying that each of these partial derivatives is continuous. Written out in more detail,

$$\mathbf{f}(\mathbf{x} + \mathbf{v}) - \mathbf{f}(\mathbf{x}) = D\mathbf{f}(\mathbf{x})\mathbf{v} + \mathbf{o}(\mathbf{v}) = \sum_{k=1}^n \frac{\partial \mathbf{f}}{\partial x_k}(\mathbf{x}) v_k + \mathbf{o}(\mathbf{v})$$

Thus $D\mathbf{f}(\mathbf{x})\mathbf{v} = \sum_{k=1}^n \frac{\partial \mathbf{f}}{\partial x_k}(\mathbf{x}) v_k$. Now consider the second derivative.

$$D^2 \mathbf{f}(\mathbf{x})(\mathbf{w})(\mathbf{v}) =$$

$$\begin{aligned} D\mathbf{f}(\mathbf{x} + \mathbf{w})\mathbf{v} - D\mathbf{f}(\mathbf{x})\mathbf{v} + \mathbf{o}(\mathbf{w})(\mathbf{v}) &= \sum_{k=1}^n \left(\frac{\partial \mathbf{f}}{\partial x_k}(\mathbf{x} + \mathbf{w}) - \frac{\partial \mathbf{f}}{\partial x_k}(\mathbf{x}) \right) v_k + \mathbf{o}(\mathbf{w})(\mathbf{v}) \\ &= \sum_{k=1}^n \left(\sum_{j=1}^n \frac{\partial^2 \mathbf{f}}{\partial x_j \partial x_k}(\mathbf{x}) w_j + \mathbf{o}(\mathbf{w}) \right) v_k + \mathbf{o}(\mathbf{w})(\mathbf{v}) = \sum_{j,k} \frac{\partial^2 \mathbf{f}}{\partial x_j \partial x_k}(\mathbf{x}) w_j v_k + \mathbf{o}(\mathbf{w})(\mathbf{v}) \end{aligned}$$

and so $D^2 \mathbf{f}(\mathbf{x})(\mathbf{w})(\mathbf{v}) = \sum_{j,k} \frac{\partial^2 \mathbf{f}}{\partial x_j \partial x_k}(\mathbf{x}) w_j v_k$. Hence $D^2 \mathbf{f}$ is continuous if and only if each of these coefficients $\mathbf{x} \rightarrow \frac{\partial^2 \mathbf{f}}{\partial x_j \partial x_k}(\mathbf{x})$ is continuous. Obviously you can continue doing this and conclude that $D^k \mathbf{f}$ is continuous if and only if all of the partial derivatives of order up to k are continuous. ■

In practice, this is usually what people are thinking when they say that \mathbf{f} is C^k . But as just argued, this is the same as saying that the r linear form $\mathbf{x} \rightarrow D^r \mathbf{f}(\mathbf{x})$ is continuous into the appropriate space of linear transformations for each $r \leq k$.

Of course the above is based on the assumption that the first k derivatives exist and gives two equivalent formulations which state that these derivatives are continuous. Can anything be said about the existence of the derivatives based on the existence and continuity of the partial derivatives? As pointed out, if the partial derivatives exist and are continuous, then the function is differentiable and has continuous derivative. However, I want to emphasize the idea of the Cartesian product.

7.9 The Derivative and the Cartesian Product

There are theorems which can be used to get differentiability of a function based on existence and continuity of the partial derivatives. A generalization of this was given above. Here a function defined on a product space is considered. It is very much like what was presented above and could be obtained as a special case but to reinforce the ideas, I will do it from scratch because certain aspects of it are important in the statement of the implicit function theorem.

The following is an important abstract generalization of the concept of partial derivative presented above. Instead of taking the derivative with respect to one variable, it is taken with respect to several but not with respect to others. This vague notion is made precise in the following definition. First here is a lemma.

Lemma 7.9.1 *Suppose U is an open set in $X \times Y$. Then the set, $U_{\mathbf{y}}$ defined by*

$$U_{\mathbf{y}} \equiv \{\mathbf{x} \in X : (\mathbf{x}, \mathbf{y}) \in U\}$$

is an open set in X . Here $X \times Y$ is a finite dimensional vector space in which the vector space operations are defined componentwise. Thus for $\mathbf{a}, \mathbf{b} \in \mathbb{F}$,

$$\mathbf{a}(\mathbf{x}_1, \mathbf{y}_1) + \mathbf{b}(\mathbf{x}_2, \mathbf{y}_2) = (\mathbf{a}\mathbf{x}_1 + \mathbf{b}\mathbf{x}_2, \mathbf{a}\mathbf{y}_1 + \mathbf{b}\mathbf{y}_2)$$

and the norm can be taken to be

$$\|(\mathbf{x}, \mathbf{y})\| \equiv \max(\|\mathbf{x}\|, \|\mathbf{y}\|)$$

Proof: Recall by Theorem 5.2.4 it does not matter how this norm is defined and the definition above is convenient. It obviously satisfies most axioms of a norm. The only one which is not obvious is the triangle inequality. I will show this now.

$$\begin{aligned} \|(\mathbf{x}, \mathbf{y}) + (\mathbf{x}_1, \mathbf{y}_1)\| &\equiv \|(\mathbf{x} + \mathbf{x}_1, \mathbf{y} + \mathbf{y}_1)\| \equiv \max(\|\mathbf{x} + \mathbf{x}_1\|, \|\mathbf{y} + \mathbf{y}_1\|) \\ &\leq \max(\|\mathbf{x}\| + \|\mathbf{x}_1\|, \|\mathbf{y}\| + \|\mathbf{y}_1\|) \\ &\leq \max(\|\mathbf{x}\|, \|\mathbf{y}\|) + \max(\|\mathbf{x}_1\|, \|\mathbf{y}_1\|) \\ &\equiv \|(\mathbf{x}, \mathbf{y})\| + \|(\mathbf{x}_1, \mathbf{y}_1)\| \end{aligned}$$

Let $\mathbf{x} \in U_{\mathbf{y}}$. Then $(\mathbf{x}, \mathbf{y}) \in U$ and so there exists $r > 0$ such that $B((\mathbf{x}, \mathbf{y}), r) \in U$. This says that if $(\mathbf{u}, \mathbf{v}) \in X \times Y$ such that $\|(\mathbf{u}, \mathbf{v}) - (\mathbf{x}, \mathbf{y})\| < r$, then $(\mathbf{u}, \mathbf{v}) \in U$. Thus if

$$\|(\mathbf{u}, \mathbf{v}) - (\mathbf{x}, \mathbf{y})\| = \|\mathbf{u} - \mathbf{x}\|_X < r,$$

then $(\mathbf{u}, \mathbf{y}) \in U$. This has just said that $B(\mathbf{x}, r)_X$, the ball taken in X is contained in $U_{\mathbf{y}}$. This proves the lemma. ■

Or course one could also consider $U_{\mathbf{x}} \equiv \{\mathbf{y} : (\mathbf{x}, \mathbf{y}) \in U\}$ in the same way and conclude this set is open in Y . Also, the generalization to many factors yields the same conclusion. In this case, for $\mathbf{x} \in \prod_{i=1}^n X_i$, let

$$\|\mathbf{x}\| \equiv \max \left(\|\mathbf{x}_i\|_{X_i} : \mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n) \right)$$

Then a similar argument to the above shows this is a norm on $\prod_{i=1}^n X_i$. Consider the triangle inequality.

$$\begin{aligned} \|(\mathbf{x}_1, \dots, \mathbf{x}_n) + (\mathbf{y}_1, \dots, \mathbf{y}_n)\| &= \max_i \left(\|\mathbf{x}_i + \mathbf{y}_i\|_{X_i} \right) \leq \max_i \left(\|\mathbf{x}_i\|_{X_i} + \|\mathbf{y}_i\|_{X_i} \right) \\ &\leq \max_i \left(\|\mathbf{x}_i\|_{X_i} \right) + \max_i \left(\|\mathbf{y}_i\|_{X_i} \right) = \|\mathbf{x}\| + \|\mathbf{y}\| \end{aligned}$$

Corollary 7.9.2 *Let $U \subseteq \prod_{i=1}^n X_i$ be an open set and let*

$$U_{(\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n)} \equiv \{\mathbf{x} \in \mathbb{R}^{r_i} : (\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n) \in U\}.$$

Then $U_{(\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n)}$ is an open set in \mathbb{R}^{r_i} .

Proof: Let $\mathbf{z} \in U_{(\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n)}$. Then $(\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{z}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n) \equiv \mathbf{x} \in U$ by definition. Therefore, since U is open, there exists $r > 0$ such that $B(\mathbf{x}, r) \subseteq U$. It follows that for $B(\mathbf{z}, r)_{X_i}$ denoting the ball in X_i , it follows that $B(\mathbf{z}, r)_{X_i} \subseteq U_{(\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n)}$ because to say that $\|\mathbf{z} - \mathbf{w}\|_{X_i} < r$ is to say that

$$\|(\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{z}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n) - (\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{w}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n)\| < r$$

and so $\mathbf{w} \in U_{(\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n)}$. ■

Next is a generalization of the partial derivative.

Definition 7.9.3 *Let $g : U \subseteq \prod_{i=1}^n X_i \rightarrow Y$, where U is an open set. Then the map*

$$\mathbf{z} \rightarrow g(\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{z}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n)$$

is a function from the open set in X_i ,

$$\{\mathbf{z} : \mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{z}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n) \in U\}$$

to Y . When this map is differentiable, its derivative is denoted by $D_i g(\mathbf{x})$. To aid in the notation, for $\mathbf{v} \in X_i$, let $\theta_i \mathbf{v} \in \prod_{i=1}^n X_i$ be the vector $(\mathbf{0}, \dots, \mathbf{v}, \dots, \mathbf{0})$ where the \mathbf{v} is in the i^{th} slot and for $\mathbf{v} \in \prod_{i=1}^n X_i$, let v_i denote the entry in the i^{th} slot of \mathbf{v} . Thus, by saying

$$\mathbf{z} \rightarrow g(\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{z}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n)$$

is differentiable is meant that for $\mathbf{v} \in X_i$ sufficiently small,

$$g(\mathbf{x} + \theta_i \mathbf{v}) - g(\mathbf{x}) = D_i g(\mathbf{x}) \mathbf{v} + o(\mathbf{v}).$$

Note $D_i g(\mathbf{x}) \in \mathcal{L}(X_i, Y)$.

As discussed above, we have the following definition of $C^1(U)$.

Definition 7.9.4 Let $U \subseteq X$ be an open set. Then $\mathbf{f} : U \rightarrow Y$ is $C^1(U)$ if \mathbf{f} is differentiable and the mapping $\mathbf{x} \rightarrow D\mathbf{f}(\mathbf{x})$, is continuous as a function from U to $\mathcal{L}(X, Y)$.

With this definition of partial derivatives, here is the major theorem. Note the resemblance with the matrix of the derivative of a function having values in \mathbb{R}^m in terms of the partial derivatives.

Theorem 7.9.5 Let $\mathbf{g}, U, \prod_{i=1}^n X_i$, be given as in Definition 7.9.3. Then \mathbf{g} is $C^1(U)$ if and only if $D_i \mathbf{g}$ exists and is continuous on U for each i . In this case, \mathbf{g} is differentiable and

$$D\mathbf{g}(\mathbf{x})(\mathbf{v}) = \sum_k D_k \mathbf{g}(\mathbf{x}) v_k \quad (7.14)$$

where $\mathbf{v} = (v_1, \dots, v_n)$.

Proof: Suppose then that $D_i \mathbf{g}$ exists and is continuous for each i . Note $\sum_{j=1}^k \theta_j v_j = (v_1, \dots, v_k, 0, \dots, 0)$. Thus $\sum_{j=1}^n \theta_j v_j = \mathbf{v}$ and define $\sum_{j=1}^0 \theta_j v_j \equiv \mathbf{0}$. Therefore,

$$\begin{aligned} \mathbf{g}(\mathbf{x} + \mathbf{v}) - \mathbf{g}(\mathbf{x}) &= \sum_{k=1}^n \left[\mathbf{g} \left(\mathbf{x} + \sum_{j=1}^k \theta_j v_j \right) - \mathbf{g} \left(\mathbf{x} + \sum_{j=1}^{k-1} \theta_j v_j \right) \right] \\ &= \sum_{k=1}^n \left[\left(\mathbf{g} \left(\mathbf{x} + \sum_{j=1}^k \theta_j v_j \right) - \mathbf{g}(\mathbf{x} + \theta_k v_k) \right) - \left(\mathbf{g} \left(\mathbf{x} + \sum_{j=1}^{k-1} \theta_j v_j \right) - \mathbf{g}(\mathbf{x}) \right) \right] \\ &\quad + \sum_{k=1}^n (\mathbf{g}(\mathbf{x} + \theta_k v_k) - \mathbf{g}(\mathbf{x})) \end{aligned} \quad (7.15)$$

If $\mathbf{h}_k(\mathbf{x}) \equiv \mathbf{g} \left(\mathbf{x} + \sum_{j=1}^{k-1} \theta_j v_j \right) - \mathbf{g}(\mathbf{x})$ then the top sum is $\sum_{k=1}^n \mathbf{h}_k(\mathbf{x} + \theta_k v_k) - \mathbf{h}_k(\mathbf{x})$ and from the definition of \mathbf{h}_k , $\|D\mathbf{h}_k(\mathbf{x})\| < \varepsilon$ a sufficiently small ball containing \mathbf{x} . Thus this top sum is dominated by $\varepsilon \|\mathbf{v}\|$ whenever $\|\mathbf{v}\|$ is small enough. Since ε is arbitrary, this term is $\mathbf{o}(\mathbf{v})$. The last term is $\sum_{k=1}^n D_k \mathbf{g}(\mathbf{x}) v_k + \mathbf{o}(v_k)$ and so, collecting terms obtains

$$\mathbf{g}(\mathbf{x} + \mathbf{v}) - \mathbf{g}(\mathbf{x}) = \sum_{k=1}^n D_k \mathbf{g}(\mathbf{x}) v_k + \mathbf{o}(\mathbf{v})$$

which shows $D\mathbf{g}(\mathbf{x})$ exists and equals the formula given in 7.14. Also $\mathbf{x} \rightarrow D\mathbf{g}(\mathbf{x})$ is continuous since each of the $D_k \mathbf{g}(\mathbf{x})$ are.

Next suppose \mathbf{g} is C^1 . I need to verify that $D_k \mathbf{g}(\mathbf{x})$ exists and is continuous. Let $\mathbf{v} \in X_k$ sufficiently small. Then

$$\mathbf{g}(\mathbf{x} + \theta_k \mathbf{v}) - \mathbf{g}(\mathbf{x}) = D\mathbf{g}(\mathbf{x}) \theta_k \mathbf{v} + \mathbf{o}(\theta_k \mathbf{v}) = D\mathbf{g}(\mathbf{x}) \theta_k \mathbf{v} + \mathbf{o}(\mathbf{v})$$

since $\|\theta_k \mathbf{v}\| = \|\mathbf{v}\|$. Then $D_k \mathbf{g}(\mathbf{x})$ exists and equals $D\mathbf{g}(\mathbf{x}) \circ \theta_k$. Now $\mathbf{x} \rightarrow D\mathbf{g}(\mathbf{x})$ is continuous. Since θ_k is linear, it follows from Lemma 5.2.1 that $\theta_k : X_k \rightarrow \prod_{i=1}^n X_i$ is also continuous. ■

Note that the above argument also works at a single point \mathbf{x} . That is, continuity at \mathbf{x} of the partials implies $D\mathbf{g}(\mathbf{x})$ exists and is continuous at \mathbf{x} .

The way this is usually used is in the following corollary which has already been obtained. Remember the matrix of $D\mathbf{f}(\mathbf{x})$. Recall that if a function is C^1 in the sense that $\mathbf{x} \rightarrow D\mathbf{f}(\mathbf{x})$ is continuous then all the partial derivatives exist and are continuous. The next corollary says that if the partial derivatives do exist and are continuous, then the function is differentiable and has continuous derivative.

Corollary 7.9.6 *Let U be an open subset of \mathbb{F}^n and let $\mathbf{f} : U \rightarrow \mathbb{F}^m$ be C^1 in the sense that all the partial derivatives of \mathbf{f} exist and are continuous. Then \mathbf{f} is differentiable and*

$$\mathbf{f}(\mathbf{x} + \mathbf{v}) = \mathbf{f}(\mathbf{x}) + \sum_{k=1}^n \frac{\partial \mathbf{f}}{\partial x_k}(\mathbf{x}) v_k + \mathbf{o}(\mathbf{v}).$$

Similarly, if the partial derivatives up to order k exist and are continuous, then the function is C^k in the sense that the first k derivatives exist and are continuous.

7.10 Mixed Partial Derivatives

Continuing with the special case where f is defined on an open set in \mathbb{F}^n , I will next consider an interesting result which was known to Euler in around 1734 about mixed partial derivatives. It was proved by Clairaut some time later. It turns out that the mixed partial derivatives, if continuous will end up being equal. Recall the notation $f_x = \frac{\partial f}{\partial x} = D_{\mathbf{e}_1}f$ and $f_{xy} = \frac{\partial^2 f}{\partial y \partial x} = D_{\mathbf{e}_1}D_{\mathbf{e}_2}f$.

Theorem 7.10.1 *Suppose $f : U \subseteq \mathbb{F}^2 \rightarrow \mathbb{R}$ where U is an open set on which f_x, f_y, f_{xy} and f_{yx} exist. Then if f_{xy} and f_{yx} are continuous at the point $(x, y) \in U$, it follows*

$$f_{xy}(x, y) = f_{yx}(x, y).$$

Proof: Since U is open, there exists $r > 0$ such that $B((x, y), r) \subseteq U$. Now let $|t|, |s| < r/2, t, s$ real numbers and consider

$$\Delta(s, t) \equiv \frac{1}{st} \left\{ \overbrace{f(x+t, y+s) - f(x+t, y)}^{h(t)} - \overbrace{(f(x, y+s) - f(x, y))}^{h(0)} \right\}. \quad (7.16)$$

Note that $(x+t, y+s) \in U$ because

$$\begin{aligned} |(x+t, y+s) - (x, y)| &= |(t, s)| = (t^2 + s^2)^{1/2} \\ &\leq \left(\frac{r^2}{4} + \frac{r^2}{4} \right)^{1/2} = \frac{r}{\sqrt{2}} < r. \end{aligned}$$

As implied above, $h(t) \equiv f(x+t, y+s) - f(x+t, y)$. Therefore, by the mean value theorem from one variable calculus and the (one variable) chain rule,

$$\begin{aligned} \Delta(s, t) &= \frac{1}{st} (h(t) - h(0)) = \frac{1}{st} h'(\alpha t) t \\ &= \frac{1}{s} (f_x(x + \alpha t, y+s) - f_x(x + \alpha t, y)) \end{aligned}$$

for some $\alpha \in (0, 1)$. Applying the mean value theorem again,

$$\Delta(s, t) = f_{xy}(x + \alpha t, y + \beta s)$$

where $\alpha, \beta \in (0, 1)$.

If the terms $f(x+t, y)$ and $f(x, y+s)$ are interchanged in 7.16, $\Delta(s, t)$ is unchanged and the above argument shows there exist $\gamma, \delta \in (0, 1)$ such that

$$\Delta(s, t) = f_{yx}(x + \gamma t, y + \delta s).$$

Letting $(s, t) \rightarrow (0, 0)$ and using the continuity of f_{xy} and f_{yx} at (x, y) ,

$$\lim_{(s,t) \rightarrow (0,0)} \Delta(s, t) = f_{xy}(x, y) = f_{yx}(x, y). \blacksquare$$

The following is obtained from the above by simply fixing all the variables except for the two of interest.

Corollary 7.10.2 Suppose U is an open subset of X and $f : U \rightarrow \mathbb{R}$ has the property that for two indices, k, l , f_{x_k} , f_{x_l} , $f_{x_l x_k}$, and $f_{x_k x_l}$ exist on U and $f_{x_k x_l}$ and $f_{x_l x_k}$ are both continuous at $\mathbf{x} \in U$. Then $f_{x_k x_l}(\mathbf{x}) = f_{x_l x_k}(\mathbf{x})$.

By considering the real and imaginary parts of f in the case where f has values in \mathbb{C} you obtain the following corollary.

Corollary 7.10.3 Suppose U is an open subset of \mathbb{F}^n and $f : U \rightarrow \mathbb{F}$ has the property that for two indices, k, l , f_{x_k} , f_{x_l} , $f_{x_l x_k}$, and $f_{x_k x_l}$ exist on U and $f_{x_k x_l}$ and $f_{x_l x_k}$ are both continuous at $\mathbf{x} \in U$. Then $f_{x_k x_l}(\mathbf{x}) = f_{x_l x_k}(\mathbf{x})$.

Finally, by considering the components of \mathbf{f} you get the following generalization.

Corollary 7.10.4 Suppose U is an open subset of \mathbb{F}^n and $\mathbf{f} : U \rightarrow \mathbb{F}^m$ has the property that for two indices, k, l , \mathbf{f}_{x_k} , \mathbf{f}_{x_l} , $\mathbf{f}_{x_l x_k}$, and $\mathbf{f}_{x_k x_l}$ exist on U and $\mathbf{f}_{x_k x_l}$ and $\mathbf{f}_{x_l x_k}$ are both continuous at $\mathbf{x} \in U$. Then $\mathbf{f}_{x_k x_l}(\mathbf{x}) = \mathbf{f}_{x_l x_k}(\mathbf{x})$.

This can be generalized to functions which have values in a normed linear space, but I plan to stop with what is given above. One way to proceed would be to reduce to a consideration of the coordinate maps and then apply the above. It would even hold in infinite dimensions through the use of the Hahn Banach theorem. The idea is to reduce to the scalar valued case as above.

In addition, it is obvious that for a function of many variables you could pick any pair and say these are equal if they are both continuous.

It is necessary to assume the mixed partial derivatives are continuous in order to assert they are equal. The following is a well known example [2].

Example 7.10.5 Let

$$f(x, y) = \begin{cases} \frac{xy(x^2 - y^2)}{x^2 + y^2} & \text{if } (x, y) \neq (0, 0) \\ 0 & \text{if } (x, y) = (0, 0) \end{cases}$$

From the definition of partial derivatives it follows that $f_x(0, 0) = f_y(0, 0) = 0$. Using the standard rules of differentiation, for $(x, y) \neq (0, 0)$,

$$f_x = y \frac{x^4 - y^4 + 4x^2 y^2}{(x^2 + y^2)^2}, \quad f_y = x \frac{x^4 - y^4 - 4x^2 y^2}{(x^2 + y^2)^2}$$

Now

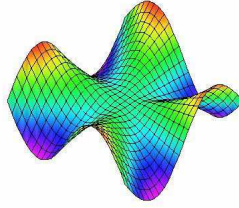
$$f_{xy}(0,0) \equiv \lim_{y \rightarrow 0} \frac{f_x(0,y) - f_x(0,0)}{y} = \lim_{y \rightarrow 0} \frac{-y^4}{(y^2)^2} = -1$$

while

$$f_{yx}(0,0) \equiv \lim_{x \rightarrow 0} \frac{f_y(x,0) - f_y(0,0)}{x} = \lim_{x \rightarrow 0} \frac{x^4}{(x^2)^2} = 1$$

showing that although the mixed partial derivatives do exist at $(0,0)$, they are not equal there.

Incidentally, the graph of this function appears very innocent. Its fundamental sickness is not apparent. It is like one of those whited sepulchers mentioned in the Bible.



7.11 A Cofactor Identity

Lemma 7.11.1 Suppose $\det(A) = 0$. Then for all sufficiently small nonzero ε , it follows that $\det(A + \varepsilon I) \neq 0$.

Proof: Let $\det(\lambda I - A) = \lambda^p + a_1\lambda^{p-1} + \cdots + a_{p-1}\lambda + a_p$. First suppose A is a $p \times p$ matrix. If $\det(A) \neq 0$, this will still be true for all ε small enough. Now suppose also that $\det(A) = 0$. Thus, the constant term of $\det(\lambda I - A)$ is 0. Consider $\varepsilon I + A \equiv A_\varepsilon$ for small real ε . The characteristic polynomial of A_ε is

$$\det(\lambda I - A_\varepsilon) = \det((\lambda - \varepsilon)I - A)$$

This is of the form

$$(\lambda - \varepsilon)^p + a_1(\lambda - \varepsilon)^{p-1} + \cdots + (\lambda - \varepsilon)^m a_m$$

where the a_j are the coefficients in the characteristic polynomial for A and $a_k = 0$ for $k > m, a_m \neq 0$. The constant term of this polynomial in λ must be nonzero for all ε small enough because it is of the form

$$(-1)^m \varepsilon^m a_m + (\text{higher order terms in } \varepsilon) = \varepsilon^m [a_m (-1)^m + \varepsilon C(\varepsilon)]$$

which is nonzero for all positive but very small ε . Thus $\varepsilon I + A$ is invertible for all ε small enough but nonzero. ■

Recall that for A an $p \times p$ matrix, $\text{cof}(A)_{ij}$ is the determinant of the matrix which results from deleting the i^{th} row and the j^{th} column and multiplying by $(-1)^{i+j}$. In the proof and in what follows, I am using $D\mathbf{g}$ to equal the matrix of the linear transformation $D\mathbf{g}$ taken with respect to the usual basis on \mathbb{R}^p . Thus $(D\mathbf{g})_{ij} = \partial g_i / \partial x_j$ where $\mathbf{g} = \sum_i g_i \mathbf{e}_i$ for the \mathbf{e}_i the standard basis vectors.

Lemma 7.11.2 Let $g : U \rightarrow \mathbb{R}^p$ be C^2 where U is an open subset of \mathbb{R}^p . Then

$$\sum_{j=1}^p \text{cof}(Dg)_{ij,j} = 0,$$

where here $(Dg)_{ij} \equiv g_{i,j} \equiv \frac{\partial g_i}{\partial x_j}$. Also, $\text{cof}(Dg)_{ij} = \frac{\partial \det(Dg)}{\partial g_{i,j}}$.

Proof: From the cofactor expansion theorem,

$$\delta_{kj} \det(Dg) = \sum_{i=1}^p g_{i,k} \text{cof}(Dg)_{ij} \quad (7.17)$$

This is because if $k \neq j$, that on the right is the cofactor expansion of a determinant with two equal columns while if $k = j$, it is just the cofactor expansion of the determinant. In particular,

$$\frac{\partial \det(Dg)}{\partial g_{i,j}} = \text{cof}(Dg)_{ij} \quad (7.18)$$

which shows the last claim of the lemma. Assume that $Dg(x)$ is invertible to begin with. Differentiate 7.17 with respect to x_j and sum on j . This yields

$$\sum_{r,s,j} \delta_{kj} \frac{\partial (\det Dg)}{\partial g_{r,s}} g_{r,s,j} = \sum_{ij} g_{i,kj} (\text{cof}(Dg))_{ij} + \sum_{ij} g_{i,k} \text{cof}(Dg)_{ij,j}.$$

Hence, using $\delta_{kj} = 0$ if $j \neq k$ and 7.18,

$$\sum_{rs} (\text{cof}(Dg))_{rs} g_{r,s,k} = \sum_{rs} g_{r,ks} (\text{cof}(Dg))_{rs} + \sum_{ij} g_{i,k} \text{cof}(Dg)_{ij,j}.$$

Subtracting the first sum on the right from both sides and using the equality of mixed partials,

$$\sum_i g_{i,k} \left(\sum_j (\text{cof}(Dg))_{ij,j} \right) = 0.$$

Since it is assumed Dg is invertible, this shows $\sum_j (\text{cof}(Dg))_{ij,j} = 0$. If $\det(Dg) = 0$, use Lemma 7.11.1 to let $g_k(x) = g(x) + \varepsilon_k x$ where $\varepsilon_k \rightarrow 0$ and $\det(Dg + \varepsilon_k I) \equiv \det(Dg_k) \neq 0$. Then

$$\sum_j (\text{cof}(Dg))_{ij,j} = \lim_{k \rightarrow \infty} \sum_j (\text{cof}(Dg_k))_{ij,j} = 0 \blacksquare$$

7.12 Newton's Method

Remember Newton's method from one variable calculus. It was an algorithm for finding the zeros of a function. Beginning with x_k the next iterate was $x_{k+1} = x_k - f'(x_k)^{-1} (f(x_k))$. Of course the same thing can sometimes work in \mathbb{R}^n or even more generally. Here you have a function $f(x)$ and you want to locate a zero. Then you could consider the sequence of iterates $x_{k+1} = x_k - Df(x_k)^{-1} (f(x_k))$. If the sequence converges to x then you would have $x = x - Df(x)^{-1} (f(x))$ and so you would need to have $f(x) = 0$. In the next section, a modification of this well known method will be used to prove the Implicit function theorem. The modification is that you look for a solution to the equation near

\mathbf{x}_0 and replace the above algorithm with the simpler one $\mathbf{x}_{k+1} = \mathbf{x}_k - D\mathbf{f}(\mathbf{x}_0)^{-1}(\mathbf{f}(\mathbf{x}_k))$. Then if $T\mathbf{x} = \mathbf{x} - D\mathbf{f}(\mathbf{x}_0)^{-1}(\mathbf{f}(\mathbf{x}))$, it follows that as long as \mathbf{x} is sufficiently close to \mathbf{x}_0 , $DT(\mathbf{x}) = I - D\mathbf{f}(\mathbf{x}_0)^{-1}D\mathbf{f}(\mathbf{x})$ and the norm of this transformation is very small so one can use the mean value inequality to conclude that T is a contraction mapping and provide a sequence of iterates which converge to a fixed point. Actually, the situation will be a little more complicated because we will do the implicit function theorem first, but this is the idea.

7.13 Exercises

1. Here are some scalar valued functions of several variables. Determine which of these functions are $o(\mathbf{v})$. Here \mathbf{v} is a vector in \mathbb{R}^n , $\mathbf{v} = (v_1, \dots, v_n)$.

- | | |
|---------------------------|---|
| (a) $v_1 v_2$ | (e) $v_1(v_1 + v_2 + xv_3)$ |
| (b) $v_2 \sin(v_1)$ | (f) $(e^{v_1} - 1 - v_1)$ |
| (c) $v_1^2 + v_2$ | (g) $(\mathbf{x} \cdot \mathbf{v}) \mathbf{v} $ |
| (d) $v_2 \sin(v_1 + v_2)$ | |

2. Here is a function of two variables. $f(x, y) = x^2 y + x^2$. Find $Df(x, y)$ directly from the definition. Recall this should be a linear transformation which results from multiplication by a 1×2 matrix. Find this matrix.

3. Let $\mathbf{f}(x, y) = \begin{pmatrix} x^2 + y \\ y^2 \end{pmatrix}$. Compute the derivative directly from the definition. This should be the linear transformation which results from multiplying by a 2×2 matrix. Find this matrix.

4. You have $\mathbf{h}(\mathbf{x}) = \mathbf{g}(\mathbf{f}(\mathbf{x}))$. Here $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^m$ and $\mathbf{g}(\mathbf{y}) \in \mathbb{R}^p$. where \mathbf{f}, \mathbf{g} are appropriately differentiable. Thus $D\mathbf{h}(\mathbf{x})$ results from multiplication by a matrix. Using the chain rule, give a formula for the ij^{th} entry of this matrix. How does this relate to multiplication of matrices? In other words, you have two matrices which correspond to $D\mathbf{g}(\mathbf{f}(\mathbf{x}))$ and $D\mathbf{f}(\mathbf{x})$. Call $\mathbf{z} = \mathbf{g}(\mathbf{y})$, $\mathbf{y} = \mathbf{f}(\mathbf{x})$. Then

$$D\mathbf{g}(\mathbf{y}) = \begin{pmatrix} \frac{\partial z}{\partial y_1} & \cdots & \frac{\partial z}{\partial y_m} \end{pmatrix}, D\mathbf{f}(\mathbf{x}) = \begin{pmatrix} \frac{\partial y}{\partial x_1} & \cdots & \frac{\partial y}{\partial x_n} \end{pmatrix}$$

Explain the manner in which the ij^{th} entry of $D\mathbf{h}(\mathbf{x})$ is $\sum_k \frac{\partial z_i}{\partial y_k} \frac{\partial y_k}{\partial x_j}$. This is a review of the way we multiply matrices. what is the i^{th} row of $D\mathbf{g}(\mathbf{y})$ and the j^{th} column of $D\mathbf{f}(\mathbf{x})$?

5. Find $f_x, f_y, f_z, f_{xy}, f_{yx}, f_{zy}$ for the following. Verify the mixed partial derivatives are equal.

- | |
|-------------------------------|
| (a) $x^2 y^3 z^4 + \sin(xyz)$ |
| (b) $\sin(xyz) + x^2 yz$ |

6. Suppose f is a continuous function and $f: U \rightarrow \mathbb{R}$ where U is an open set and suppose that $\mathbf{x} \in U$ has the property that for all \mathbf{y} near \mathbf{x} , $f(\mathbf{x}) \leq f(\mathbf{y})$. Prove that if f has all of its partial derivatives at \mathbf{x} , then $f_{x_i}(\mathbf{x}) = 0$ for each x_i . **Hint:** Consider $f(\mathbf{x} + t\mathbf{v}) = h(t)$. Argue that $h'(0) = 0$ and then see what this implies about $Df(\mathbf{x})$.

7. As an important application of Problem 6 consider the following. Experiments are done at n times, t_1, t_2, \dots, t_n and at each time there results a collection of numerical outcomes. Denote by $\{(t_i, x_i)\}_{i=1}^p$ the set of all such pairs and try to find numbers a and b such that the line $x = at + b$ approximates these ordered pairs as well as possible in the sense that out of all choices of a and b , $\sum_{i=1}^p (at_i + b - x_i)^2$ is as small as possible. In other words, you want to minimize the function of two variables $f(a, b) \equiv \sum_{i=1}^p (at_i + b - x_i)^2$. Find a formula for a and b in terms of the given ordered pairs. You will be finding the formula for the least squares regression line.
8. Let f be a function which has continuous derivatives. Show that $u(t, x) = f(x - ct)$ solves the wave equation $u_{tt} - c^2 \Delta u = 0$. What about $u(x, t) = f(x + ct)$? Here $\Delta u = u_{xx}$.
9. Show that if $\Delta u = \lambda u$ where u is a function of only x , then $e^{\lambda t} u$ solves the heat equation $u_t - \Delta u = 0$. Here $\Delta u = u_{xx}$.
10. Show that if $f(x) = o(x)$, then $f'(0) = 0$.
11. Let $f(x, y)$ be defined on \mathbb{R}^2 as follows. $f(x, x^2) = 1$ if $x \neq 0$. Define $f(0, 0) = 0$, and $f(x, y) = 0$ if $y \neq x^2$. Show that f is not continuous at $(0, 0)$ but that

$$\lim_{h \rightarrow 0} \frac{f(ha, hb) - f(0, 0)}{h} = 0$$

for (a, b) an arbitrary vector. Thus the Gateaux derivative exists at $(0, 0)$ in every direction but f is not even continuous there.

12. Let

$$f(x, y) \equiv \begin{cases} \frac{xy^4}{x^2 + y^8} & \text{if } (x, y) \neq (0, 0) \\ 0 & \text{if } (x, y) = (0, 0) \end{cases}$$

Show that this function is not continuous at $(0, 0)$ but that the Gateaux derivative $\lim_{h \rightarrow 0} \frac{f(ha, hb) - f(0, 0)}{h}$ exists and equals 0 for every vector (a, b) .

13. Let U be an open subset of \mathbb{R}^n and suppose that $f : [a, b] \times U \rightarrow \mathbb{R}$ satisfies

$$(x, y) \rightarrow \frac{\partial f}{\partial y_i}(x, y), (x, y) \rightarrow f(x, y)$$

are all continuous. Show that $\int_a^b f(x, y) dx$, $\int_a^b \frac{\partial f}{\partial y_i}(x, y) dx$ all make sense and that in fact

$$\frac{\partial}{\partial y_i} \left(\int_a^b f(x, y) dx \right) = \int_a^b \frac{\partial f}{\partial y_i}(x, y) dx$$

Also explain why $y \rightarrow \int_a^b \frac{\partial f}{\partial y_i}(x, y) dx$ is continuous. **Hint:** You will need to use the theorems from one variable calculus about the existence of the integral for a continuous function. You may also want to use theorems about uniform continuity of continuous functions defined on compact sets.

14. I found this problem in Apostol's book [1]. This is a very important result and is obtained very simply. Read it and fill in any missing details. Let $g(x) \equiv \int_0^1 \frac{e^{-x^2(1+t^2)}}{1+t^2} dt$ and $f(x) \equiv \left(\int_0^x e^{-t^2} dt \right)^2$. Note $\frac{\partial}{\partial x} \left(\frac{e^{-x^2(1+t^2)}}{1+t^2} \right) = -2xe^{-x^2(1+t^2)}$. Explain why this is so. Also show the conditions of Problem 13 are satisfied so that

$$g'(x) = \int_0^1 \left(-2xe^{-x^2(1+t^2)} \right) dt.$$

Now use the chain rule and the fundamental theorem of calculus to find $f'(x)$. Then change the variable in the formula for $f'(x)$ to make it an integral from 0 to 1 and show $f'(x) + g'(x) = 0$. Now this shows $f(x) + g(x)$ is a constant. Show the constant is $\pi/4$ by letting $x \rightarrow 0$. Next take a limit as $x \rightarrow \infty$ to obtain the following formula for the improper integral, $\int_0^\infty e^{-t^2} dt, \left(\int_0^\infty e^{-t^2} dt \right)^2 = \pi/4$. In passing to the limit in the integral for g as $x \rightarrow \infty$ you need to justify why that integral converges to 0. To do this, argue the integrand converges uniformly to 0 for $t \in [0, 1]$ and then explain why this gives convergence of the integral. Thus $\int_0^\infty e^{-t^2} dt = \sqrt{\pi}/2$.

15. Recall the treatment of integrals of continuous functions in Proposition 5.9.5 or what you used in beginning calculus. The gamma function is defined for $x > 0$ as $\Gamma(x) \equiv \int_0^\infty e^{-t} t^{x-1} dt \equiv \lim_{R \rightarrow \infty} \int_0^R e^{-t} t^{x-1} dt$. Show this limit exists. Note you might have to give a meaning to $\int_0^R e^{-t} t^{x-1} dt$ if $x < 1$. Also show that $\Gamma(x+1) = x\Gamma(x)$, $\Gamma(1) = 1$. How does $\Gamma(n)$ for n an integer compare with $(n-1)!$?
16. Show the mean value theorem for integrals. Suppose $f \in C([a, b])$. Then there exists $x \in (a, b)$, not just in $[a, b]$ such that $f(x)(b-a) = \int_a^b f(t) dt$. **Hint:** Let $F(x) \equiv \int_a^x f(t) dt$ and use the mean value theorem, Theorem 5.9.3 along with $F'(x) = f(x)$.
17. Show, using the Weierstrass approximation theorem that linear combinations of the form $\sum_{i,j} a_{ij} g_i(s) h_j(t)$ where g_i, h_j are continuous functions on $[0, b]$ are dense in $C([0, b] \times [0, b])$, the continuous functions defined on $[0, b] \times [0, b]$ with norm given by

$$\|f\| \equiv \max \{ |f(x, y)| : (x, y) \in [0, b] \times [0, b] \}$$

Show that for h, g continuous, $\int_0^b \int_0^s g(s) h(t) dt ds - \int_0^b \int_t^b g(s) h(t) ds dt = 0$. Now explain why if f is in $C([0, b] \times [0, b])$,

$$\int_0^b \int_0^s f(s, t) dt ds - \int_0^b \int_t^b f(s, t) ds dt = 0.$$

18. Let $f(x) \equiv \left(\int_0^x e^{-t^2} dt \right)^2$. Use Proposition 5.9.5 which includes the fundamental theorem of calculus and elementary change of variables, show that

$$f'(x) = 2e^{-x^2} \left(\int_0^x e^{-t^2} dt \right) = 2e^{-x^2} \left(\int_0^1 e^{-(xs)^2} x ds \right) = \int_0^1 2xe^{-x^2(1+s^2)} ds.$$

Now show

$$f(x) = \int_0^1 \int_0^x 2te^{-t^2(1+s^2)} dt ds.$$

Show $\lim_{x \rightarrow \infty} \int_0^x e^{-t^2} dt = \frac{1}{2}\sqrt{\pi}$

Chapter 8

Implicit Function Theorem

8.1 Statement and Proof of the Theorem

Recall the following notation. $\mathcal{L}(X, Y)$ is the space of bounded linear mappings from X to Y where here $(X, \|\cdot\|_X)$ and $(Y, \|\cdot\|_Y)$ are normed linear spaces. Recall that this means that for each $L \in \mathcal{L}(X, Y)$, $\|L\| \equiv \sup_{\|x\| \leq 1} \|Lx\| < \infty$. As shown earlier, this makes $\mathcal{L}(X, Y)$ into a normed linear space. In case X is finite dimensional, $\mathcal{L}(X, Y)$ is the same as the collection of linear maps from X to Y . This was shown earlier. In what follows X, Y will be Banach spaces. If you like, think of them as finite dimensional normed linear spaces, but if you like more generality, just think: complete normed linear space and $\mathcal{L}(X, Y)$ is the space of **bounded** linear maps. In either case, this symbol is given in the following definition.

Definition 8.1.1 Let $(X, \|\cdot\|_X)$ and $(Y, \|\cdot\|_Y)$ be two normed linear spaces. Then $\mathcal{L}(X, Y)$ denotes the set of linear maps from X to Y which also satisfy the following condition. For $L \in \mathcal{L}(X, Y)$,

$$\lim_{\|x\|_X \leq 1} \|Lx\|_Y \equiv \|L\| < \infty$$

Recall that this operator norm is less than infinity is always the case where X is finite dimensional. However, if you wish to consider infinite dimensional situations, you assume the operator norm is finite as a qualification for being in $\mathcal{L}(X, Y)$.

Definition 8.1.2 A complete normed linear space is called a Banach space.

Theorem 8.1.3 If Y is a Banach space, then $\mathcal{L}(X, Y)$ is also a Banach space.

Proof: Let $\{L_n\}$ be a Cauchy sequence in $\mathcal{L}(X, Y)$ and let $x \in X$.

$$\|L_n x - L_m x\| \leq \|x\| \|L_n - L_m\|.$$

Thus $\{L_n x\}$ is a Cauchy sequence. Let $Lx = \lim_{n \rightarrow \infty} L_n x$. Then, clearly, L is linear because if x_1, x_2 are in X , and a, b are scalars, then

$$\begin{aligned} L(ax_1 + bx_2) &= \lim_{n \rightarrow \infty} L_n(ax_1 + bx_2) = \lim_{n \rightarrow \infty} (aL_n x_1 + bL_n x_2) \\ &= aLx_1 + bLx_2. \end{aligned}$$

Also L is bounded. To see this, note that $\{\|L_n\|\}$ is a Cauchy sequence of real numbers because $|\|L_n\| - \|L_m\|| \leq \|L_n - L_m\|$. Hence there exists $K > \sup\{\|L_n\| : n \in \mathbb{N}\}$. Thus, if $x \in X$, $\|Lx\| = \lim_{n \rightarrow \infty} \|L_n x\| \leq K\|x\|$. ■

The following theorem is really nice. The series in this theorem is called the Neuman series.

Lemma 8.1.4 Let $(X, \|\cdot\|)$ is a Banach space, and if $A \in \mathcal{L}(X, X)$ and $\|A\| = r < 1$, then $(I - A)^{-1} = \sum_{k=0}^{\infty} A^k \in \mathcal{L}(X, X)$ where the series converges in the Banach space $\mathcal{L}(X, X)$. If O consists of the invertible maps in $\mathcal{L}(X, X)$, then O is open and if \mathfrak{I} is the mapping which takes A to A^{-1} , then \mathfrak{I} is continuous.

Proof: First of all, why does the series make sense?

$$\left\| \sum_{k=p}^q A^k \right\| \leq \sum_{k=p}^q \|A^k\| \leq \sum_{k=p}^q \|A\|^k \leq \sum_{k=p}^{\infty} r^k \leq \frac{r^p}{1-r}$$

and so the partial sums are Cauchy in $\mathcal{L}(X, X)$. Therefore, the series converges to something in $\mathcal{L}(X, X)$ by completeness of this normed linear space. Now why is it the inverse?

$$\begin{aligned} \sum_{k=0}^{\infty} A^k (I - A) &\equiv \lim_{n \rightarrow \infty} \sum_{k=0}^n A^k (I - A) = \lim_{n \rightarrow \infty} \left(\sum_{k=0}^n A^k - \sum_{k=1}^{n+1} A^k \right) \\ &= \lim_{n \rightarrow \infty} (I - A^{n+1}) = I \end{aligned}$$

because $\|A^{n+1}\| \leq \|A\|^{n+1} \leq r^{n+1}$. Similarly,

$$(I - A) \sum_{k=0}^{\infty} A^k = \lim_{n \rightarrow \infty} (I - A^{n+1}) = I$$

and so this shows that this series is indeed the desired inverse.

Next suppose $A \in O$ so $A^{-1} \in \mathcal{L}(X, X)$. Then suppose $\|A - B\| < \frac{r}{1 + \|A^{-1}\|}$, $r < 1$. Does it follow that B is also invertible? $B = A - (A - B) = A [I - A^{-1}(A - B)]$. Then $\|A^{-1}(A - B)\| \leq \|A^{-1}\| \|A - B\| < r$ and so $[I - A^{-1}(A - B)]^{-1}$ exists. Hence $B^{-1} = [I - A^{-1}(A - B)]^{-1} A^{-1}$. Thus O is open as claimed. As to continuity, let A, B be as just described. Then using the Neuman series,

$$\begin{aligned} \|\mathfrak{I}A - \mathfrak{I}B\| &= \left\| A^{-1} - [I - A^{-1}(A - B)]^{-1} A^{-1} \right\| \\ &= \left\| A^{-1} - \sum_{k=0}^{\infty} (A^{-1}(A - B))^k A^{-1} \right\| = \left\| \sum_{k=1}^{\infty} (A^{-1}(A - B))^k A^{-1} \right\| \\ &\leq \sum_{k=1}^{\infty} \|A^{-1}\|^{k+1} \|A - B\|^k = \|A - B\| \|A^{-1}\|^2 \sum_{k=0}^{\infty} \|A^{-1}\|^k \left(\frac{r}{1 + \|A^{-1}\|} \right)^k \\ &\leq \|B - A\| \|A^{-1}\|^2 \frac{1}{1 - r}. \end{aligned}$$

Thus \mathfrak{I} is continuous at $A \in O$. ■

Next features the inverse in which there are two different spaces.

Lemma 8.1.5 *Let*

$$O \equiv \{A \in \mathcal{L}(X, Y) : A^{-1} \in \mathcal{L}(Y, X)\}$$

and let $\mathfrak{I} : O \rightarrow \mathcal{L}(Y, X)$, $\mathfrak{I}A \equiv A^{-1}$. Then O is open and \mathfrak{I} is in $C^m(O)$ for all $m = 1, 2, \dots$. Also

$$D\mathfrak{I}(A)(B) = -\mathfrak{I}(A)(B)\mathfrak{I}(A). \quad (8.1)$$

In particular, \mathfrak{I} is continuous.

Proof: Let $A \in O$ and let $B \in \mathcal{L}(X, Y)$ with $\|B\| \leq \frac{1}{2} \|A^{-1}\|^{-1}$. Then

$$\|A^{-1}B\| \leq \|A^{-1}\| \|B\| \leq \frac{1}{2}$$

So by Lemma 8.1.4,

$$\begin{aligned} (A+B)^{-1} &= (I+A^{-1}B)^{-1} A^{-1} = \sum_{n=0}^{\infty} (-1)^n (A^{-1}B)^n A^{-1} \\ &= [I - A^{-1}B + o(B)] A^{-1} \end{aligned}$$

which shows that O is open and, also,

$$\begin{aligned} \mathfrak{J}(A+B) - \mathfrak{J}(A) &= \sum_{n=0}^{\infty} (-1)^n (A^{-1}B)^n A^{-1} - A^{-1} \\ &= -A^{-1}BA^{-1} + o(B) \\ &= -\mathfrak{J}(A)(B)\mathfrak{J}(A) + o(B) \end{aligned}$$

which demonstrates 8.1. The reason the left over material is $o(B)$ follows from the observation that $o(B)$ is $\sum_{n=2}^{\infty} (-1)^n (A^{-1}B)^n A^{-1}$ and so

$$\left\| \sum_{n=2}^{\infty} (-1)^n (A^{-1}B)^n A^{-1} \right\| \leq \sum_{n=2}^{\infty} \left\| (A^{-1}B)^n A^{-1} \right\| \leq \|A^{-1}\| \|A^{-1}\|^2 \|B\|^2 \sum_{n=0}^{\infty} \frac{1}{2^n}$$

It follows from this that we can continue taking derivatives of \mathfrak{J} . For $\|B_1\|$ small,

$$\begin{aligned} & -[D\mathfrak{J}(A+B_1)(B) - D\mathfrak{J}(A)(B)] = \\ & \mathfrak{J}(A+B_1)(B)\mathfrak{J}(A+B_1) - \mathfrak{J}(A)(B)\mathfrak{J}(A) \\ &= \mathfrak{J}(A+B_1)(B)\mathfrak{J}(A+B_1) - \mathfrak{J}(A)(B)\mathfrak{J}(A+B_1) + \\ & \quad \mathfrak{J}(A)(B)\mathfrak{J}(A+B_1) - \mathfrak{J}(A)(B)\mathfrak{J}(A) \\ &= [\mathfrak{J}(A)(B_1)\mathfrak{J}(A) + o(B_1)](B)\mathfrak{J}(A+B_1) + \\ & \quad \mathfrak{J}(A)(B)[\mathfrak{J}(A)(B_1)\mathfrak{J}(A) + o(B_1)] \\ &= [\mathfrak{J}(A)(B_1)\mathfrak{J}(A) + o(B_1)](B)[A^{-1} - A^{-1}B_1A^{-1} + o(B_1)] + \\ & \quad \mathfrak{J}(A)(B)[\mathfrak{J}(A)(B_1)\mathfrak{J}(A) + o(B_1)] \\ &= \mathfrak{J}(A)(B_1)\mathfrak{J}(A)(B)\mathfrak{J}(A) + \mathfrak{J}(A)(B)\mathfrak{J}(A)(B_1)\mathfrak{J}(A) + o(B_1) \end{aligned}$$

and so

$$D^2\mathfrak{J}(A)(B_1)(B) = \mathfrak{J}(A)(B_1)\mathfrak{J}(A)(B)\mathfrak{J}(A) + \mathfrak{J}(A)(B)\mathfrak{J}(A)(B_1)\mathfrak{J}(A)$$

which shows \mathfrak{J} is $C^2(O)$. Clearly we can continue in this way which shows \mathfrak{J} is in $C^m(O)$ for all $m = 1, 2, \dots$. ■

Here are the two fundamental results presented earlier which will make it easy to prove the implicit function theorem. First is the fundamental mean value inequality.

Theorem 8.1.6 Suppose U is an open subset of X and $\mathbf{f} : U \rightarrow Y$ is differentiable on U and $\mathbf{x} + t(\mathbf{y} - \mathbf{x}) \in U$ for all $t \in [0, 1]$. (The line segment joining the two points lies in U .) Suppose also that for all points on this line segment,

$$\|D\mathbf{f}(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))\| \leq M.$$

Then

$$\|\mathbf{f}(\mathbf{y}) - \mathbf{f}(\mathbf{x})\| \leq M|\mathbf{y} - \mathbf{x}|.$$

Next recall the following theorem about fixed points of a contraction map. It was Corollary 3.8.3.

Corollary 8.1.7 Let B be a closed subset of the complete metric space (X, d) and let $f : B \rightarrow X$ be a contraction map

$$d(f(x), f(\hat{x})) \leq rd(x, \hat{x}), \quad r < 1.$$

Also suppose **there exists** $x_0 \in B$ such that the sequence of iterates $\{f^n(x_0)\}_{n=1}^\infty$ remains in B . Then f has a unique fixed point in B which is the limit of the sequence of iterates. This is a point $x \in B$ such that $f(x) = x$. In the case that $B = \overline{B(x_0, \delta)}$, the sequence of iterates satisfies the inequality

$$d(f^n(x_0), x_0) \leq \frac{d(x_0, f(x_0))}{1 - r}$$

and so it will remain in B if

$$\frac{d(x_0, f(x_0))}{1 - r} < \delta.$$

The implicit function theorem deals with the question of solving, $\mathbf{f}(\mathbf{x}, \mathbf{y}) = \mathbf{0}$ for \mathbf{x} in terms of \mathbf{y} and how smooth the solution is. It is one of the most important theorems in mathematics. The proof I will give holds with no change in the context of infinite dimensional complete normed vector spaces when suitable modifications are made on what is meant by $\mathcal{L}(X, Y)$. There are also even more general versions of this theorem than to normed vector spaces.

Recall that for X, Y normed vector spaces, the norm on $X \times Y$ is of the form

$$\|(\mathbf{x}, \mathbf{y})\| = \max(\|\mathbf{x}\|, \|\mathbf{y}\|).$$

Theorem 8.1.8 (implicit function theorem) Let X, Y, Z be finite dimensional normed vector spaces and suppose U is an open set in $X \times Y$. Let $\mathbf{f} : U \rightarrow Z$ be in $C^1(U)$ and suppose

$$\mathbf{f}(\mathbf{x}_0, \mathbf{y}_0) = \mathbf{0}, \quad D_1\mathbf{f}(\mathbf{x}_0, \mathbf{y}_0)^{-1} \in \mathcal{L}(Z, X). \quad (8.2)$$

Then there exist positive constants, δ, η , such that for every $\mathbf{y} \in B(\mathbf{y}_0, \eta)$ there exists a unique $\mathbf{x}(\mathbf{y}) \in B(\mathbf{x}_0, \delta)$ such that

$$\mathbf{f}(\mathbf{x}(\mathbf{y}), \mathbf{y}) = \mathbf{0}. \quad (8.3)$$

Furthermore, the mapping, $\mathbf{y} \rightarrow \mathbf{x}(\mathbf{y})$ is in $C^1(B(\mathbf{y}_0, \eta))$.

Proof: Let $T(x, y) \equiv x - D_1 f(x_0, y_0)^{-1} f(x, y)$. Therefore, $T(x_0, y_0) = x_0$ and

$$D_1 T(x, y) = I - D_1 f(x_0, y_0)^{-1} D_1 f(x, y). \quad (8.4)$$

by continuity of the derivative which implies continuity of $D_1 T$, it follows there exists $\delta > 0$ such that if $\|x - x_0\| < \delta$ and $\|y - y_0\| < \delta$, then

$$\|D_1 T(x, y)\| < \frac{1}{2}, \quad D_1 f(x, y)^{-1} \text{ exists} \quad (8.5)$$

The second claim follows from Lemma 8.1.5. By the mean value inequality, Theorem 8.1.6, whenever $x, x' \in B(x_0, \delta)$ and $y \in B(y_0, \delta)$,

$$\|T(x, y) - T(x', y)\| \leq \frac{1}{2} \|x - x'\|. \quad (8.6)$$

Also, it can be assumed δ is small enough that for some M and all such (x, y) ,

$$\left\| D_1 f(x_0, y_0)^{-1} \right\| \|D_2 f(x, y)\| < M \quad (8.7)$$

Next, consider only y such that $\|y - y_0\| < \eta$ where $\eta < \delta$ is so small that

$$\|T(x_0, y) - x_0\| < \frac{\delta}{3}$$

Then for such y , consider the mapping $T_y(x) = T(x, y)$. Thus by Corollary 8.1.7, for each $n \in \mathbb{N}$,

$$\delta > \frac{2}{3} \delta \geq \frac{\|T_y(x_0) - x_0\|}{1 - (1/2)} \geq \|T_y^n(x_0) - x_0\|$$

Then by 8.6, the sequence of iterations of this map T_y converges to a unique fixed point $x(y)$ in the ball $B(x_0, \delta)$. Thus, from the definition of T , $f(x(y), y) = 0$. This is the implicitly defined function.

Next we show that this function is Lipschitz continuous. For y, \hat{y} in $B(y_0, \eta)$,

$$\begin{aligned} \|T(x, y) - T(x, \hat{y})\| &= \left\| D_1 f(x_0, y_0)^{-1} f(x, y) - D_1 f(x_0, y_0)^{-1} f(x, \hat{y}) \right\| \\ &\leq M \|y - \hat{y}\| \end{aligned}$$

thanks to the above estimate 8.7 and the mean value inequality, Theorem 8.1.6. Note how convexity of $B(y_0, \eta)$ which says that the line segment joining y, \hat{y} is contained in $B(y_0, \eta)$ is important to use this theorem. Then from this,

$$\begin{aligned} \|x(y) - x(\hat{y})\| &= \|T(x(y), y) - T(x(\hat{y}), \hat{y})\| \leq \|T(x(y), y) - T(x(y), \hat{y})\| \\ &\quad + \|T(x(y), \hat{y}) - T(x(\hat{y}), \hat{y})\| \\ &\leq M \|y - \hat{y}\| + \frac{1}{2} \|x(y) - x(\hat{y})\| \end{aligned}$$

Hence,

$$\|x(y) - x(\hat{y})\| \leq 2M \|y - \hat{y}\| \quad (8.8)$$

Finally consider the claim that this implicitly defined function is C^1 .

$$\begin{aligned} 0 &= f(x(y+u), y+u) - f(x(y), y) \\ &= D_1 f(x(y), y)(x(y+u) - x(y)) + D_2 f(x(y), y)u \\ &\quad + o(x(y+u) - x(y), u) \end{aligned} \tag{8.9}$$

Consider the last term. $o(x(y+u) - x(y), u) / \|u\|$ equals

$$\begin{cases} \frac{o(x(y+u) - x(y), u)}{\|x(y+u) - x(y), u\|_{X \times Y}} \frac{\max(\|x(y+u) - x(y)\|, \|u\|)}{\|u\|} & \text{if } \|x(y+u) - x(y), u\|_{X \times Y} \neq 0 \\ 0 & \text{if } \|x(y+u) - x(y), u\|_{X \times Y} = 0 \end{cases}$$

Now the Lipschitz condition just established shows that

$$\frac{\max(\|x(y+u) - x(y)\|, \|u\|)}{\|u\|}$$

is bounded for nonzero u sufficiently small that $y, y+u \in B(y_0, \eta)$. Therefore,

$$\lim_{u \rightarrow 0} \frac{o(x(y+u) - x(y), u)}{\|u\|} = 0$$

Then 8.9 shows that

$$0 = D_1 f(x(y), y)(x(y+u) - x(y)) + D_2 f(x(y), y)u + o(u)$$

Therefore, solving for $x(y+u) - x(y)$, it follows that

$$\begin{aligned} x(y+u) - x(y) &= -D_1 f(x(y), y)^{-1} D_2 f(x(y), y)u + D_1 f(x(y), y)^{-1} o(u) \\ &= -D_1 f(x(y), y)^{-1} D_2 f(x(y), y)u + o(u) \end{aligned}$$

and now, the continuity of the partial derivatives $D_1 f, D_2 f$, continuity of the map $A \rightarrow A^{-1}$, along with the continuity of $y \rightarrow x(y)$ shows that $y \rightarrow x(y)$ is C^1 with derivative equal to $-D_1 f(x(y), y)^{-1} D_2 f(x(y), y)$. ■

The following is a nice result on functional dependence which is an application of the implicit function theorem. See Widder [59].

Example 8.1.9 Suppose f, g are C^1 near $(x_0, y_0) \in \mathbb{R}^2$ and

Suppose f, g are C^1 and

1. $\det \begin{pmatrix} f_x(x, y) & f_y(x, y) \\ g_x(x, y) & g_y(x, y) \end{pmatrix} \neq 0$ near (x_0, y_0)
2. $f_x(x_0, y_0) \neq 0$.

Then there is a C^1 function F such that $g(x, y) = F(f(x, y))$ for (x, y) near (x_0, y_0) .

Consider $f(x, y) - z = 0$ where $z_0 \equiv f(x_0, y_0)$. By assumption and implicit function theorem, there is a C^1 function $(y, z) \rightarrow \phi(y, z)$ so that for (y, z) near (y_0, z_0) the x which solves $f(x, y) - z = 0$ is $x = \phi(y, z)$. In particular, for (x, y) close to (x_0, y_0) , $f(\phi(y, f(x, y)), y) - f(x, y) = 0$ and so

$$\phi(y, f(x, y)) = x. \tag{*}$$

Also, for (y, z) near (y_0, z_0) , it follows $f_x(\phi(y, z), y) \neq 0$.

Since $f(\phi(y, f(x, y)), y) - z = 0$, $f_x(\phi(y, z), y)\phi_y(y, z) + f_y(\phi(y, z), y) = 0$. It follows from 1. that

$$\begin{aligned} \frac{\partial}{\partial y} g(\phi(y, z), y) &= g_x(\phi(y, z), y)\phi_y(y, z) + g_y(\phi(y, z), y) \\ &= -g_x(\phi(y, z), y) \left(\frac{f_y(\phi(y, z), y)}{f_x(\phi(y, z), y)} \right) + g_y(\phi(y, z), y) \\ &= \frac{1}{f_x(\phi(y, z), y)} \begin{pmatrix} -g_x(\phi(y, z), y)f_y(\phi(y, z), y) \\ +g_y(\phi(y, z), y)f_x(\phi(y, z), y) \end{pmatrix} = 0 \end{aligned}$$

Therefore, $g(\phi(y, z), y)$ does not depend on y near $(y_0, z_0) = (y_0, f(x_0, y_0))$. Thus there exists a C^1 function $z \rightarrow F(z)$ for z near $f(x_0, y_0)$ such that $g(\phi(y, z), y) = F(z)$. From *, for (x, y) near (x_0, y_0) ,

$$g(x, y) = g(\phi(y, f(x, y)), y) = F(f(x, y))$$

Note that if $g(x, y) = F(f(x, y))$, then $(g_x, g_y) = F'(f(x, y))(f_x, f_y)$ and so the above determinant will equal 0.

The next theorem is a very important special case of the implicit function theorem known as the inverse function theorem. Actually one can also obtain the implicit function theorem from the inverse function theorem. It is done this way in [36], [39] and in [2].

Theorem 8.1.10 (inverse function theorem) *Let $x_0 \in U$, an open set in X , and let $f : U \rightarrow Y$ where X, Y are finite dimensional normed vector spaces. Suppose*

$$f \text{ is } C^1(U), \text{ and } Df(x_0)^{-1} \in \mathcal{L}(Y, X). \quad (8.10)$$

Then there exist open sets W , and V such that

$$x_0 \in W \subseteq U, \quad (8.11)$$

$$f : W \rightarrow V \text{ is one to one and onto,} \quad (8.12)$$

$$f^{-1} \text{ is } C^1, \quad (8.13)$$

Proof: Apply the implicit function theorem to the function $F(x, y) \equiv f(x) - y$ where $y_0 \equiv f(x_0)$. Thus the function $y \rightarrow x(y)$ defined in that theorem is f^{-1} . Now let $W \equiv B(x_0, \delta) \cap f^{-1}(B(y_0, \eta))$ and $V \equiv B(y_0, \eta)$. This proves the theorem. ■

8.2 More Derivatives

When you consider a C^k function f defined on an open set U , you obtain the following

$$Df(x) \in \mathcal{L}(X, Y), D^2f(x) \in \mathcal{L}(X, \mathcal{L}(X, Y)), D^3f(x) \in \mathcal{L}(X, \mathcal{L}(X, \mathcal{L}(X, Y)))$$

and so forth. Thus they can each be considered as a linear transformation with values in some vector space. When you consider the vector spaces, you see that these can also be considered as multilinear functions on X with values in Y . Now consider the product of two linear transformations $A(y)B(y)w$, where everything is given to make sense and here w

is an appropriate vector. Then if each of these linear transformations can be differentiated, you would do the following simple computation.

$$\begin{aligned}
 & (A(\mathbf{y} + \mathbf{u})B(\mathbf{y} + \mathbf{u}) - A(\mathbf{y})B(\mathbf{y}))(\mathbf{w}) \\
 &= (A(\mathbf{y} + \mathbf{u})B(\mathbf{y} + \mathbf{u}) - A(\mathbf{y})B(\mathbf{y} + \mathbf{u}) + A(\mathbf{y})B(\mathbf{y} + \mathbf{u}) - A(\mathbf{y})B(\mathbf{y}))(\mathbf{w}) \\
 &= ((DA(\mathbf{y})\mathbf{u} + \mathbf{o}(\mathbf{u}))B(\mathbf{y} + \mathbf{u}) + A(\mathbf{y})(DB(\mathbf{y})\mathbf{u} + \mathbf{o}(\mathbf{u}))) (\mathbf{w}) \\
 &= (DA(\mathbf{y})(\mathbf{u})B(\mathbf{y} + \mathbf{u}) + A(\mathbf{y})DB(\mathbf{y})(\mathbf{u}) + \mathbf{o}(\mathbf{u})) (\mathbf{w}) \\
 &= (DA(\mathbf{y})(\mathbf{u})B(\mathbf{y}) + A(\mathbf{y})DB(\mathbf{y})(\mathbf{u}) + \mathbf{o}(\mathbf{u})) (\mathbf{w})
 \end{aligned}$$

Then

$$\mathbf{u} \rightarrow (DA(\mathbf{y})(\mathbf{u})B(\mathbf{y}) + A(\mathbf{y})DB(\mathbf{y})(\mathbf{u})) (\mathbf{w})$$

is clearly linear and

$$(\mathbf{u}, \mathbf{w}) \rightarrow (DA(\mathbf{y})(\mathbf{u})B(\mathbf{y}) + A(\mathbf{y})DB(\mathbf{y})(\mathbf{u})) (\mathbf{w})$$

is bilinear and continuous as a function of \mathbf{y} . By this we mean that for a fixed choice of (\mathbf{u}, \mathbf{w}) the resulting Y valued function just described is continuous. Now if each of AB, DA, DB can be differentiated, you could replace \mathbf{y} with $\mathbf{y} + \hat{\mathbf{u}}$ and do a similar computation to obtain as many differentiations as desired, the k^{th} differentiation yielding a k linear function. You can do this as long as A and B have derivatives. Now in the case of the implicit function theorem, you have

$$Dx(\mathbf{y}) = -D_1\mathbf{f}(x(\mathbf{y}), \mathbf{y})^{-1} D_2\mathbf{f}(x(\mathbf{y}), \mathbf{y}). \quad (8.14)$$

By Lemma 8.1.5 and the implicit function theorem and the chain rule, this is the situation just discussed. Thus $D^2x(\mathbf{y})$ can be obtained. Then the formula for it will only involve Dx which is known to be continuous. Thus one can continue in this way finding derivatives till \mathbf{f} fails to have them. The inverse map never creates difficulties because it is differentiable of order m for any m thanks to Lemma 8.1.5. Thus one can conclude the following corollary.

Corollary 8.2.1 *In the implicit and inverse function theorems, you can replace C^1 with C^k in the statements of the theorems for any $k \in \mathbb{N}$.*

8.3 The Case of \mathbb{R}^n

In many applications of the implicit function theorem,

$$\mathbf{f} : U \subseteq \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$$

and $\mathbf{f}(x_0, \mathbf{y}_0) = \mathbf{0}$ while \mathbf{f} is C^1 . How can you recognize the condition of the implicit function theorem which says $D_1\mathbf{f}(x_0, \mathbf{y}_0)^{-1}$ exists? This is really not hard. You recall the matrix of the transformation $D_1\mathbf{f}(x_0, \mathbf{y}_0)$ with respect to the usual basis vectors is

$$\begin{pmatrix} f_{1,x_1}(x_0, \mathbf{y}_0) & \cdots & f_{1,x_n}(x_0, \mathbf{y}_0) \\ \vdots & & \vdots \\ f_{n,x_1}(x_0, \mathbf{y}_0) & \cdots & f_{n,x_n}(x_0, \mathbf{y}_0) \end{pmatrix}$$

and so $D_1\mathbf{f}(x_0, \mathbf{y}_0)^{-1}$ exists exactly when the determinant of the above matrix is nonzero. This is the condition to check. In the general case, you just need to verify $D_1\mathbf{f}(x_0, \mathbf{y}_0)$ is one to one and this can also be accomplished by looking at the matrix of the transformation with respect to some bases on X and Z .

8.4 Exercises

1. Let $A \in \mathcal{L}(X, Y)$. Let $f(x) \equiv Ax$. Verify from the definition that $Df(x) = A$. What if $f(x) = y + Ax$? Note the similarity with functions of a single variable.
2. You have a level surface given by

$$f(x, y, z) = 0, \quad f \text{ is } C^1(U), (x, y, z) \in U,$$

The question is whether this deserves to be called a surface. Using the implicit function theorem, show that if $f(x_0, y_0, z_0) = 0$ and if $\frac{\partial f}{\partial z}(x_0, y_0, z_0) \neq 0$ then in some open subset of \mathbb{R}^3 , the relation $f(x, y, z) = 0$ can be “solved” for z getting say $z = z(x, y)$ such that $f(x, y, z(x, y)) = 0$. What happens if $\frac{\partial f}{\partial x}(x_0, y_0, z_0) \neq 0$ or $\frac{\partial f}{\partial y}(x_0, y_0, z_0) \neq 0$? Explain why z is a C^1 map for (x, y) in some open set.

3. Let $\mathbf{x}(t) = (x(t), y(t), z(t))^T$ be a vector valued function defined for $t \in (a, b)$. Then $D\mathbf{x}(t) \in \mathcal{L}(\mathbb{R}, \mathbb{R}^3)$. We usually denote this simply as $\mathbf{x}'(t)$. Thus, considered as a matrix, it is the 3×1 matrix $(x'(t), y'(t), z'(t))^T$ the T indicating that you take the transpose. Don't worry too much about this. You can also consider this as a vector. What is the geometric significance of this vector? The answer is that this vector is tangent to the curve traced out by $\mathbf{x}(t)$ for $t \in (a, b)$. Explain why this is so using the definition of the derivative. You need to describe what is meant by being tangent first. By saying that the line $\mathbf{x} = \mathbf{a} + t\mathbf{b}$ is tangent to a parametric curve consisting of points traced out by $\mathbf{x}(t)$ for $t \in (-\delta, \delta)$ at the point $\mathbf{a} = \mathbf{x}(t)$ which is on both the line and the curve, you would want to have

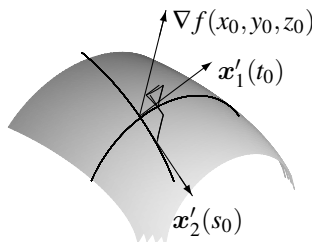
$$\lim_{u \rightarrow 0} \frac{\mathbf{x}(t+u) - (\mathbf{a} + u\mathbf{b})}{u} = \mathbf{0}$$

With this definition of what it means for a line to be tangent, explain why the line $\mathbf{x}(t) + \mathbf{x}'(t)u$ for $u \in (-\delta, \delta)$ is tangent to the curve determined by $t \rightarrow \mathbf{x}(t)$ at the point $\mathbf{x}(t)$. So why would you take the above as a definition of what it means to be tangent? Consider the component functions of $\mathbf{x}(t)$. What does the above limit say about the component functions and the corresponding components of \mathbf{b} in terms of slopes of lines tangent to curves?

4. Let $f(x, y, z)$ be a C^1 function $f : U \rightarrow \mathbb{R}$ where U is an open set in \mathbb{R}^3 . The gradient vector, defined as

$$\left(\frac{\partial f}{\partial x}(x, y, z) \quad \frac{\partial f}{\partial y}(x, y, z) \quad \frac{\partial f}{\partial z}(x, y, z) \right)^T$$

has fundamental geometric significance illustrated by the following picture.



The way we present this in engineering math is to consider a smooth C^1 curve $(x(t), y(t), z(t))$ for $t \in (a, b)$ such that when $t = c \in (a, b)$, $(x(c), y(c), z(c))$ equals the point (x, y, z) in the level surface and such that $(x(t), y(t), z(t))$ lies in this surface. Then $0 = f(x(t), y(t), z(t))$. Show, using the chain rule, that the gradient vector at the point (x, y, z) is perpendicular to

$$(x'(c), y'(c), z'(c)).$$

Recall that the chain rule says that for $h(t) = f(x(t), y(t), z(t))$, $Dh(t) =$

$$\left(\frac{\partial f}{\partial x}(x(t), y(t), z(t)) \quad \frac{\partial f}{\partial y}(x(t), y(t), z(t)) \quad \frac{\partial f}{\partial z}(x(t), y(t), z(t)) \right) \begin{pmatrix} x' \\ y' \\ z' \end{pmatrix}$$

Since this holds for all such smooth curves in the surface which go through the given point, we say that the gradient vector is perpendicular to the surface. In the picture, there are two intersecting curves which are shown to intersect at a point of the surface. We present this to the students in engineering math and everyone is happy with it, but the argument is specious. Why do there exist any smooth curves in the surface through a point? What would you need to assume to justify the existence of smooth curves in the surface at some point of the level surface? Why?

5. This problem illustrates what can happen when the gradient of a scalar valued function vanishes or is not well defined. Consider the level surface given by

$$z - \sqrt{(x^2 + y^2)} = 0.$$

Sketch the graph of this surface. Why is there no unique tangent plane at the origin $(0, 0, 0)$? Next consider $z^2 - (x^2 + y^2) = 0$. What about a well defined tangent plane at $(0, 0, 0)$?

6. Suppose you have two level surfaces $f(x, y, z) = 0$ and $g(x, y, z) = 0$ which intersect at a point (x_0, y_0, z_0) , each f, g is C^1 . Use the implicit function theorem to give conditions which will guarantee that the intersection of these two surfaces near this point is a curve. Explain why.
7. Let X, Y be Banach spaces and let U be an open subset of X . Let $f: U \rightarrow Y$ be $C^1(U)$, let $x_0 \in U$, and $\delta > 0$ be given. Show there exists $\varepsilon > 0$ such that if $x_1, x_2 \in B(x_0, \varepsilon)$, then

$$\|f(x_1) - (f(x_2) + Df(x_0)(x_1 - x_2))\| \leq \delta \|x_1 - x_2\|$$

Hint: You know $f(x_1) - f(x_2) = Df(x_2)(x_1 - x_2) + o(x_1 - x_2)$. Use continuity.

8. ↑ This problem illustrates how if $Df(x_0)$ is one to one, then near x_0 the same is true of f . Suppose in this problem that all normed linear spaces are finite dimensional. Suppose $Df(x_0)$ is one to one. Here $f: U \rightarrow Y$ where $U \subseteq X$.
- (a) Show that there exists $r > 0$ such that $\|Df(x_0)x\| \geq r\|x\|$. To do this, recall equivalence of norms.
- (b) Use the above problem to show that there is $\varepsilon > 0$ such that f is one to one on $B(x_0, \varepsilon)$ provided $Df(x_0)$ is one to one.

9. If U, V are open sets in Banach spaces X, Y respectively and $f : U \rightarrow V$ is one to one and onto and both f, f^{-1} are C^1 , show that $Df(x) : X \rightarrow Y$ is one to one and onto for each $x \in U$. **Hint:** $f \circ f^{-1} = \text{identity}$. Now use chain rule.
10. A function $f : U \subseteq \mathbb{C} \rightarrow \mathbb{C}$ where U is an open set subset of the complex numbers \mathbb{C} is called analytic if

$$\lim_{h \rightarrow 0} \frac{f(z+h) - f(z)}{h} \equiv f'(z), \quad z = x + iy$$

exists and $z \rightarrow f'(z)$ is continuous. Show that if f is analytic on an open set U and if $f'(z) \neq 0$, then there is an open set V containing z such that $f(V)$ is open, f is one to one, and f, f^{-1} are both continuous. **Hint:** This follows very easily from the inverse function theorem. Recall that we have allowed for the field of scalars the complex numbers.

11. Problem 8 has to do with concluding that f is locally one to one if $Df(x_0)$ is only known to be one to one. The next obvious question concerns the situation where $Df(x_0)$ maybe is possibly not one to one but is onto. There are two parts, a linear algebra consideration, followed by an application of the inverse function theorem. Thus these two problems are each generalizations of the inverse function theorem.
- (a) Suppose X is a finite dimensional vector space and $M \in \mathcal{L}(X, Y)$ is onto Y . Consider a basis for $M(X) = Y, \{Mx_1, \dots, Mx_n\}$. Verify that $\{x_1, \dots, x_n\}$ is linearly independent. Define $\hat{X} \equiv \text{span}(x_1, \dots, x_n)$. Show that if \hat{M} is the restriction of M to \hat{X} , then \hat{M} is one to one and onto Y .
- (b) Now suppose $f : U \subseteq X \rightarrow Y$ is C^1 and $Df(x_0)$ is onto Y . Show that there is a ball $B(f(x_0), \varepsilon)$ and an open set $V \subseteq X$ such that $f(V) \supseteq B(f(x_0), \varepsilon)$ so that if $Df(x)$ is onto for each $x \in U$, then $f(U)$ is an open set. This is called the open map theorem. You might use the inverse function theorem with the spaces \hat{X}, Y . You might want to consider Problem 1. This is a nice illustration of why we developed the inverse and implicit function theorems on arbitrary normed linear spaces. You will see that this is a fairly easy problem.
12. Recall that a function $f : U \subseteq X \rightarrow Y$ where here assume X is finite dimensional, is Gateaux differentiable if

$$\lim_{t \rightarrow 0} \frac{f(x+tv) - f(x)}{t} \equiv D_v f(x)$$

exists. Here $t \in \mathbb{R}$. Suppose that $x \rightarrow D_v f(x)$ exists and is continuous on U . Show it follows that f is differentiable and in fact $D_v f(x) = Df(x)v$. **Hint:** Let $g(y) \equiv f(\sum_i y_i x_i)$ and argue that the partial derivatives of g all exist and are continuous. Conclude that g is C^1 and then argue that f is just the composition of g with a linear map.

13. Let

$$f(x, y) = \begin{cases} \frac{(x^2 - y^4)^2}{(x^2 + y^4)^2} & \text{if } (x, y) \neq (0, 0) \\ 1 & \text{if } (x, y) = (0, 0) \end{cases}$$

Show that f is not differentiable, and in fact is not even continuous, but $D_v f(0, 0)$ exists and equals 0 for every $v \neq 0$.

14. Let

$$f(x, y) = \begin{cases} \frac{xy^4}{x^2+y^8} & \text{if } (x, y) \neq (0, 0) \\ 0 & \text{if } (x, y) = (0, 0) \end{cases}$$

Show that f is not differentiable, and in fact is not even continuous, but $D_v f(0, 0)$ exists and equals 0 for every $v \neq 0$.

8.5 The Method of Lagrange Multipliers

As an application of the implicit function theorem, consider the method of Lagrange multipliers from calculus. Recall the problem is to maximize or minimize a function subject to equality constraints. Let $f : U \rightarrow \mathbb{R}$ be a C^1 function where $U \subseteq \mathbb{R}^n$ and let

$$g_i(\mathbf{x}) = 0, \quad i = 1, \dots, m \quad (8.15)$$

be a collection of equality constraints with $m < n$. Now consider the system of nonlinear equations

$$\begin{aligned} f(\mathbf{x}) &= a \\ g_i(\mathbf{x}) &= 0, \quad i = 1, \dots, m. \end{aligned}$$

\mathbf{x}_0 is a local maximum if $f(\mathbf{x}_0) \geq f(\mathbf{x})$ for all \mathbf{x} near \mathbf{x}_0 which also satisfies the constraints 8.15. A local minimum is defined similarly. Let $\mathbf{F} : U \times \mathbb{R} \rightarrow \mathbb{R}^{m+1}$ be defined by

$$\mathbf{F}(\mathbf{x}, a) \equiv \begin{pmatrix} f(\mathbf{x}) - a \\ g_1(\mathbf{x}) \\ \vdots \\ g_m(\mathbf{x}) \end{pmatrix}. \quad (8.16)$$

Now consider the $m+1 \times n$ Jacobian matrix, the matrix of the linear transformation, $D_1 \mathbf{F}(\mathbf{x}, a)$ with respect to the usual basis for \mathbb{R}^n and \mathbb{R}^{m+1} .

$$\begin{pmatrix} f_{x_1}(\mathbf{x}_0) & \cdots & f_{x_n}(\mathbf{x}_0) \\ g_{1x_1}(\mathbf{x}_0) & \cdots & g_{1x_n}(\mathbf{x}_0) \\ \vdots & & \vdots \\ g_{mx_1}(\mathbf{x}_0) & \cdots & g_{mx_n}(\mathbf{x}_0) \end{pmatrix}.$$

If this matrix has rank $m+1$ then some $m+1 \times m+1$ submatrix has nonzero determinant. It follows from the implicit function theorem that there exist $m+1$ variables, $x_{i_1}, \dots, x_{i_{m+1}}$ such that the system

$$\mathbf{F}(\mathbf{x}, a) = \mathbf{0} \quad (8.17)$$

specifies these $m+1$ variables as a function of the remaining $n - (m+1)$ variables and a in an open set of \mathbb{R}^{n-m} . Thus there is a solution (\mathbf{x}, a) to 8.17 for some \mathbf{x} close to \mathbf{x}_0 whenever a is in some open interval. Therefore, \mathbf{x}_0 cannot be either a local minimum or a local maximum. It follows that if \mathbf{x}_0 is either a local maximum or a local minimum, then the above matrix must have rank less than $m+1$ which requires the rows to be linearly dependent. Thus, there exist m scalars,

$$\lambda_1, \dots, \lambda_m,$$

and a scalar μ , not all zero such that

$$\mu \begin{pmatrix} f_{x_1}(\mathbf{x}_0) \\ \vdots \\ f_{x_n}(\mathbf{x}_0) \end{pmatrix} = \lambda_1 \begin{pmatrix} g_{1x_1}(\mathbf{x}_0) \\ \vdots \\ g_{1x_n}(\mathbf{x}_0) \end{pmatrix} + \cdots + \lambda_m \begin{pmatrix} g_{mx_1}(\mathbf{x}_0) \\ \vdots \\ g_{mx_n}(\mathbf{x}_0) \end{pmatrix}. \quad (8.18)$$

If the column vectors

$$\begin{pmatrix} g_{1x_1}(\mathbf{x}_0) \\ \vdots \\ g_{1x_n}(\mathbf{x}_0) \end{pmatrix}, \cdots, \begin{pmatrix} g_{mx_1}(\mathbf{x}_0) \\ \vdots \\ g_{mx_n}(\mathbf{x}_0) \end{pmatrix} \quad (8.19)$$

are linearly independent, then, $\mu \neq 0$ and dividing by μ yields an expression of the form

$$\begin{pmatrix} f_{x_1}(\mathbf{x}_0) \\ \vdots \\ f_{x_n}(\mathbf{x}_0) \end{pmatrix} = \lambda_1 \begin{pmatrix} g_{1x_1}(\mathbf{x}_0) \\ \vdots \\ g_{1x_n}(\mathbf{x}_0) \end{pmatrix} + \cdots + \lambda_m \begin{pmatrix} g_{mx_1}(\mathbf{x}_0) \\ \vdots \\ g_{mx_n}(\mathbf{x}_0) \end{pmatrix} \quad (8.20)$$

at every point \mathbf{x}_0 which is either a local maximum or a local minimum. This proves the following theorem.

Theorem 8.5.1 *Let U be an open subset of \mathbb{R}^n and let $f : U \rightarrow \mathbb{R}$ be a C^1 function. Then if $\mathbf{x}_0 \in U$ is either a local maximum or local minimum of f subject to the constraints 8.15, then 8.18 must hold for some scalars $\mu, \lambda_1, \dots, \lambda_m$ not all equal to zero. If the vectors in 8.19 are linearly independent, it follows that an equation of the form 8.20 holds.*

8.6 The Taylor Formula

First recall the Taylor formula with the Lagrange form of the remainder. It will only be needed on $[0, 1]$ so that is what I will show.

Theorem 8.6.1 *Let $h : [0, 1] \rightarrow \mathbb{R}$ have $m+1$ derivatives. Then there exists $t \in (0, 1)$ such that*

$$h(1) = h(0) + \sum_{k=1}^m \frac{h^{(k)}(0)}{k!} + \frac{h^{(m+1)}(t)}{(m+1)!}.$$

Proof: Let K be a number chosen such that

$$h(1) - \left(h(0) + \sum_{k=1}^m \frac{h^{(k)}(0)}{k!} + K \right) = 0$$

Now the idea is to find K . To do this, let

$$F(t) = h(1) - \left(h(t) + \sum_{k=1}^m \frac{h^{(k)}(t)}{k!} (1-t)^k + K(1-t)^{m+1} \right)$$

Then $F(1) = F(0) = 0$. Therefore, by Rolle's theorem or the mean value theorem, Theorem 5.9.3, there exists t between 0 and 1 such that $F'(t) = 0$. Thus,

$$\begin{aligned} 0 &= -F'(t) = h'(t) + \sum_{k=1}^m \frac{h^{(k+1)}(t)}{k!} (1-t)^k \\ &\quad - \sum_{k=1}^m \frac{h^{(k)}(t)}{k!} k(1-t)^{k-1} - K(m+1)(1-t)^m \end{aligned}$$

And so

$$\begin{aligned}
 &= h'(t) + \sum_{k=1}^m \frac{h^{(k+1)}(t)}{k!} (1-t)^k - \sum_{k=0}^{m-1} \frac{h^{(k+1)}(t)}{k!} (1-t)^k \\
 &\quad - K(m+1)(1-t)^m \\
 &= h'(t) + \frac{h^{(m+1)}(t)}{m!} (1-t)^m - h'(t) - K(m+1)(1-t)^m
 \end{aligned}$$

and so $K = \frac{h^{(m+1)}(t)}{(m+1)!}$. This proves the theorem. ■

Now let $f : U \rightarrow \mathbb{R}$ where $U \subseteq X$ a normed vector space and suppose $f \in C^m(U)$ and suppose $D^{m+1}f(x)$ exists on U . Let $x \in U$ and let $r > 0$ be such that

$$B(x, r) \subseteq U.$$

Then for $\|v\| < r$ consider

$$f(x+tv) - f(x) \equiv h(t)$$

for $t \in [0, 1]$. Then by the chain rule,

$$h'(t) = Df(x+tv)(v); \quad h''(t) = D^2f(x+tv)(v)(v)$$

and continuing in this way,

$$h^{(k)}(t) = D^{(k)}f(x+tv)(v)(v) \cdots (v) \equiv D^{(k)}f(x+tv)v^k.$$

It follows from Taylor's formula for a function of one variable given above that

$$f(x+v) = f(x) + \sum_{k=1}^m \frac{D^{(k)}f(x)v^k}{k!} + \frac{D^{(m+1)}f(x+tv)v^{m+1}}{(m+1)!}. \quad (8.21)$$

This proves the following theorem.

Theorem 8.6.2 *Let $f : U \rightarrow \mathbb{R}$ and let $f \in C^{m+1}(U)$. Then if*

$$B(x, r) \subseteq U,$$

and $\|v\| < r$, there exists $t \in (0, 1)$ such that 8.21 holds.

8.7 Second Derivative Test

Now consider the case where $U \subseteq \mathbb{R}^n$ and $f : U \rightarrow \mathbb{R}$ is $C^2(U)$. Then from Taylor's theorem, if v is small enough, there exists $t \in (0, 1)$ such that

$$f(x+v) = f(x) + Df(x)v + \frac{D^2f(x+tv)v^2}{2}. \quad (8.22)$$

Consider

$$\begin{aligned}
 D^2f(x+tv)(e_i)(e_j) &\equiv D(Df(x+tv))e_i e_j \\
 &= D\left(\frac{\partial f(x+tv)}{\partial x_i}\right)e_j = \frac{\partial^2 f(x+tv)}{\partial x_j \partial x_i}
 \end{aligned}$$

where e_i are the usual basis vectors. Letting $\mathbf{v} = \sum_{i=1}^n v_i e_i$, the second derivative term in 8.22 reduces to

$$\frac{1}{2} \sum_{i,j} D^2 f(\mathbf{x} + t\mathbf{v})(e_i)(e_j) v_i v_j = \frac{1}{2} \sum_{i,j} H_{ij}(\mathbf{x} + t\mathbf{v}) v_i v_j$$

where

$$H_{ij}(\mathbf{x} + t\mathbf{v}) = D^2 f(\mathbf{x} + t\mathbf{v})(e_i)(e_j) = \frac{\partial^2 f(\mathbf{x} + t\mathbf{v})}{\partial x_j \partial x_i}.$$

Definition 8.7.1 The matrix whose ij^{th} entry is $\frac{\partial^2 f(\mathbf{x})}{\partial x_j \partial x_i}$ is called the Hessian matrix, denoted as $\mathbf{H}(\mathbf{x})$.

From Theorem 7.10.1, this is a symmetric real matrix, thus self adjoint. By the continuity of the second partial derivative,

$$\begin{aligned} f(\mathbf{x} + \mathbf{v}) &= f(\mathbf{x}) + Df(\mathbf{x})\mathbf{v} + \frac{1}{2}\mathbf{v}^T \mathbf{H}(\mathbf{x})\mathbf{v} + \\ &\quad \frac{1}{2}(\mathbf{v}^T (\mathbf{H}(\mathbf{x} + t\mathbf{v}) - \mathbf{H}(\mathbf{x}))\mathbf{v}). \end{aligned} \quad (8.23)$$

where the last two terms involve ordinary matrix multiplication and

$$\mathbf{v}^T = \begin{pmatrix} v_1 & \cdots & v_n \end{pmatrix}$$

for v_i the components of \mathbf{v} relative to the standard basis.

Definition 8.7.2 Let $f : D \rightarrow \mathbb{R}$ where D is a subset of some normed vector space. Then f has a local minimum at $\mathbf{x} \in D$ if there exists $\delta > 0$ such that for all $\mathbf{y} \in B(\mathbf{x}, \delta)$

$$f(\mathbf{y}) \geq f(\mathbf{x}).$$

f has a local maximum at $\mathbf{x} \in D$ if there exists $\delta > 0$ such that for all $\mathbf{y} \in B(\mathbf{x}, \delta)$

$$f(\mathbf{y}) \leq f(\mathbf{x}).$$

Theorem 8.7.3 If $f : U \rightarrow \mathbb{R}$ where U is an open subset of \mathbb{R}^n and f is C^2 , suppose $Df(\mathbf{x}) = 0$. Then if $\mathbf{H}(\mathbf{x})$ has all positive eigenvalues, \mathbf{x} is a local minimum. If the Hessian matrix $\mathbf{H}(\mathbf{x})$ has all negative eigenvalues, then \mathbf{x} is a local maximum. If $\mathbf{H}(\mathbf{x})$ has a positive eigenvalue, then there exists a direction in which f has a local minimum at \mathbf{x} , while if $\mathbf{H}(\mathbf{x})$ has a negative eigenvalue, there exists a direction in which $\mathbf{H}(\mathbf{x})$ has a local maximum at \mathbf{x} .

Proof: Since $Df(\mathbf{x}) = 0$, formula 8.23 holds and by continuity of the second derivative, $\mathbf{H}(\mathbf{x})$ is a symmetric matrix. Thus $\mathbf{H}(\mathbf{x})$ has all real eigenvalues. Suppose first that $\mathbf{H}(\mathbf{x})$ has all positive eigenvalues and that all are larger than $\delta^2 > 0$. Then by Theorem 1.4.1, $\mathbf{H}(\mathbf{x})$ has an orthonormal basis of eigenvectors, $\{v_i\}_{i=1}^n$ and if \mathbf{u} is an arbitrary vector, such that $\mathbf{u} = \sum_{j=1}^n u_j v_j$ where $u_j = \mathbf{u} \cdot v_j$, then

$$\mathbf{u}^T \mathbf{H}(\mathbf{x}) \mathbf{u} = \sum_{j=1}^n u_j v_j^T \mathbf{H}(\mathbf{x}) \sum_{j=1}^n u_j v_j$$

$$= \sum_{j=1}^n u_j^2 \lambda_j \geq \delta^2 \sum_{j=1}^n u_j^2 = \delta^2 |u|^2.$$

From 8.23 and the continuity of H , if v is small enough,

$$f(x+v) \geq f(x) + \frac{1}{2} \delta^2 |v|^2 - \frac{1}{4} \delta^2 |v|^2 = f(x) + \frac{\delta^2}{4} |v|^2.$$

This shows the first claim of the theorem. The second claim follows from similar reasoning. Suppose $H(x)$ has a positive eigenvalue λ^2 . Then let v be an eigenvector for this eigenvalue. Then from 8.23,

$$\begin{aligned} f(x+tv) &= f(x) + \frac{1}{2} t^2 v^T H(x) v + \\ &\quad \frac{1}{2} t^2 (v^T (H(x+tv) - H(x)) v) \end{aligned}$$

which implies

$$\begin{aligned} f(x+tv) &= f(x) + \frac{1}{2} t^2 \lambda^2 |v|^2 + \frac{1}{2} t^2 (v^T (H(x+tv) - H(x)) v) \\ &\geq f(x) + \frac{1}{4} t^2 \lambda^2 |v|^2 \end{aligned}$$

whenever t is small enough. Thus in the direction v the function has a local minimum at x . The assertion about the local maximum in some direction follows similarly. This proves the theorem. ■

This theorem is an analogue of the second derivative test for higher dimensions. As in one dimension, when there is a zero eigenvalue, it may be impossible to determine from the Hessian matrix what the local qualitative behavior of the function is. For example, consider

$$f_1(x, y) = x^4 + y^2, \quad f_2(x, y) = -x^4 + y^2.$$

Then $Df_i(0, 0) = \mathbf{0}$ and for both functions, the Hessian matrix evaluated at $(0, 0)$ equals

$$\begin{pmatrix} 0 & 0 \\ 0 & 2 \end{pmatrix}$$

but the behavior of the two functions is very different near the origin. The second has a saddle point while the first has a minimum there.

8.8 The Rank Theorem

This is a very interesting result. The proof follows Marsden and Hoffman. First here is some linear algebra.

Theorem 8.8.1 *Let $L: \mathbb{R}^n \rightarrow \mathbb{R}^N$ have rank m . Then there exists a basis*

$$\{u_1, \dots, u_m, u_{m+1}, \dots, u_n\}$$

such that a basis for $\ker(L)$ is $\{u_{m+1}, \dots, u_n\}$.

Proof: Since L has rank m , there is a basis for $L(\mathbb{R}^n)$ which is of the form

$$\{Lu_1, \dots, Lu_m\}$$

Then if $\sum_i c_i u_i = 0$ you can do L to both sides and conclude that each $c_i = 0$. Hence $\{u_1, \dots, u_m\}$ is linearly independent. Let $\{v_1, \dots, v_k\}$ be a basis for $\ker(L)$. Let $x \in \mathbb{R}^n$. Then $Lx = \sum_{i=1}^m c_i Lu_i$ for some choice of scalars c_i . Hence $L(x - \sum_{i=1}^m c_i u_i) = 0$ which shows that there exist d_j such that $x = \sum_{i=1}^m c_i u_i + \sum_{j=1}^k d_j v_j$. It follows that

$$\text{span}(u_1, \dots, u_m, v_1, \dots, v_k) = \mathbb{R}^n$$

Is this set of vectors linearly independent? Suppose $\sum_{i=1}^m c_i u_i + \sum_{j=1}^k d_j v_j = 0$. Do L to both sides to get $\sum_{i=1}^m c_i Lu_i = 0$. Thus each $c_i = 0$. Hence $\sum_{j=1}^k d_j v_j = 0$ and so each $d_j = 0$ also. It follows that $k = n - m$ and we can let

$$\{v_1, \dots, v_k\} = \{u_{m+1}, \dots, u_n\}. \blacksquare$$

Another useful linear algebra result is the following lemma.

Lemma 8.8.2 *Let $V \subseteq \mathbb{R}^n$ be a subspace and suppose $A(x) \in \mathcal{L}(V, \mathbb{R}^N)$ for x in some open set U . Also suppose $x \rightarrow A(x)$ is continuous for $x \in U$. Then if $A(x_0)$ is one to one on V for some $x_0 \in U$, then it follows that for all x close enough to x_0 , $A(x)$ is also one to one on V .*

Proof: Consider V as an inner product space with the inner product from \mathbb{R}^n and $A(x)^* A(x)$. Then $A(x)^* A(x) \in \mathcal{L}(V, V)$ and $x \rightarrow A(x)^* A(x)$ is also continuous. Also for $v \in V$,

$$(A(x)^* A(x)v, v)_V = (A(x)v, A(x)v)_{\mathbb{R}^N}$$

If $A(x_0)^* A(x_0)v = 0$, then from the above, it follows that $A(x_0)v = 0$ also. Therefore, $v = 0$ and so $A(x_0)^* A(x_0)$ is one to one on V . For all x close enough to x_0 , it follows from continuity that $A(x)^* A(x)$ is also one to one. Thus, for such x , if $A(x)v = 0$, Then $A(x)^* A(x)v = 0$ and so $v = 0$. Thus, for x close enough to x_0 , it follows that $A(x)$ is also one to one on V . \blacksquare

Theorem 8.8.3 *Let $f: A \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^N$ where A is open in \mathbb{R}^n . Let f be a C^r function and suppose that $Df(x)$ has rank m for all $x \in A$. Let $x_0 \in A$. Then there are open sets $U, V \subseteq \mathbb{R}^n$ with $x_0 \in V$, and a C^r function $h: U \rightarrow V$ with inverse $h^{-1}: V \rightarrow U$ also C^r such that $f \circ h$ depends only on (x_1, \dots, x_m) .*

Proof: Let $L = Df(x_0)$, and $N_0 = \ker L$. Using the above linear algebra theorem, there exists

$$\{u_1, \dots, u_m\}$$

such that $\{Lu_1, \dots, Lu_m\}$ is a basis for $L\mathbb{R}^n$. Extend to form a basis for \mathbb{R}^n ,

$$\{u_1, \dots, u_m, u_{m+1}, \dots, u_n\}$$

such that a basis for $N_0 = \ker L$ is $\{u_{m+1}, \dots, u_n\}$. Let

$$M \equiv \text{span}(u_1, \dots, u_m).$$

Let the coordinate maps be ψ_k so that if $\mathbf{x} \in \mathbb{R}^n$,

$$\mathbf{x} = \psi_1(\mathbf{x})\mathbf{u}_1 + \cdots + \psi_n(\mathbf{x})\mathbf{u}_n$$

Since these coordinate maps are linear, they are infinitely differentiable.

Next I will define coordinate maps for $\mathbf{x} \in \mathbb{R}^N$. Then by the above construction, $\{L\mathbf{u}_1, \dots, L\mathbf{u}_m\}$ is a basis for $L(\mathbb{R}^n)$. Let a basis for \mathbb{R}^N be

$$\{L\mathbf{u}_1, \dots, L\mathbf{u}_m, \mathbf{v}_{m+1}, \dots, \mathbf{v}_N\}$$

(Note that, since the rank of $D\mathbf{f}(\mathbf{x}) = m$ you must have $N \geq m$.) The coordinate maps ϕ_i will be defined as follows for $\mathbf{x} \in \mathbb{R}^N$.

$$\mathbf{x} = \phi_1(\mathbf{x})L\mathbf{u}_1 + \cdots + \phi_m(\mathbf{x})L\mathbf{u}_m + \phi_{m+1}(\mathbf{x})\mathbf{v}_{m+1} + \cdots + \phi_N(\mathbf{x})\mathbf{v}_N$$

Now define two infinitely differentiable maps $G: \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $H: \mathbb{R}^N \rightarrow \mathbb{R}^n$,

$$G(\mathbf{x}) \equiv (0, \dots, 0, \psi_{m+1}(\mathbf{x}), \dots, \psi_n(\mathbf{x}))$$

$$H(\mathbf{y}) \equiv (\phi_1(\mathbf{y}), \dots, \phi_m(\mathbf{y}), 0, \dots, 0)$$

For $\mathbf{x} \in A \subseteq \mathbb{R}^n$, let

$$\mathbf{g}(\mathbf{x}) \equiv H(\mathbf{f}(\mathbf{x})) + G(\mathbf{x}) \in \mathbb{R}^n$$

Thus the first term picks out the first m entries of $\mathbf{f}(\mathbf{x})$ and the second term the last $n - m$ entries of \mathbf{x} . It is of the form

$$(\phi_1(\mathbf{f}(\mathbf{x})), \dots, \phi_m(\mathbf{f}(\mathbf{x})), \psi_{m+1}(\mathbf{x}), \dots, \psi_n(\mathbf{x}))$$

Then

$$D\mathbf{g}(\mathbf{x}_0)(\mathbf{v}) = H\mathbf{L}(\mathbf{v}) + G\mathbf{v} = H\mathbf{L}\mathbf{v} + G\mathbf{v} \quad (8.24)$$

which is of the form

$$D\mathbf{g}(\mathbf{x}_0)(\mathbf{v}) = (\phi_1(L\mathbf{v}), \dots, \phi_m(L\mathbf{v}), \psi_{m+1}(\mathbf{v}), \dots, \psi_n(\mathbf{v}))$$

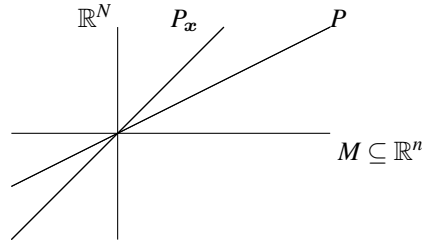
If this equals $\mathbf{0}$, then all the components of \mathbf{v} , $\psi_{m+1}(\mathbf{v}), \dots, \psi_n(\mathbf{v})$ are equal to 0. Hence

$$\mathbf{v} = \sum_{i=1}^m c_i \mathbf{u}_i.$$

But also the coordinates of $L\mathbf{v}, \phi_1(L\mathbf{v}), \dots, \phi_m(L\mathbf{v})$ are all zero so $L\mathbf{v} = \mathbf{0}$ and so $\mathbf{0} = \sum_{i=1}^m c_i L\mathbf{u}_i$ so by independence of the $L\mathbf{u}_i$, each $c_i = 0$ and consequently $\mathbf{v} = \mathbf{0}$.

This proves the conditions for the inverse function theorem are valid for \mathbf{g} . Therefore, there is an open ball U and an open set V , $\mathbf{x}_0 \in V$, such that $\mathbf{g}: V \rightarrow U$ is a C^r map and its inverse $\mathbf{g}^{-1}: U \rightarrow V$ is also. We can assume by continuity and Lemma 8.8.2 that V and U are small enough that for each $\mathbf{x} \in V$, $D\mathbf{g}(\mathbf{x})$ is one to one. This follows from the fact that $\mathbf{x} \rightarrow D\mathbf{g}(\mathbf{x})$ is continuous.

Since it is assumed that $D\mathbf{f}(\mathbf{x})$ is of rank m , $D\mathbf{f}(\mathbf{x})(\mathbb{R}^n)$ is a subspace which is m dimensional, denoted as $P_{\mathbf{x}}$. Also denote $L(\mathbb{R}^n) = L(M)$ as P .



Thus $\{L\mathbf{u}_1, \dots, L\mathbf{u}_m\}$ is a basis for P . Using Lemma 8.8.2 again, by making V, U smaller if necessary, one can also assume that for each $\mathbf{x} \in V$, $D\mathbf{f}(\mathbf{x})$ is one to one on M (although not on \mathbb{R}^n) and $HD\mathbf{f}(\mathbf{x})$ is one to one on M . This follows from continuity and the fact that $L = D\mathbf{f}(\mathbf{x}_0)$ is one to one on M . Therefore, it is also the case that $D\mathbf{f}(\mathbf{x})$ maps the m dimensional space M onto the m dimensional space $P_{\mathbf{x}}$ and H is one to one on $P_{\mathbf{x}}$. The reason for this last claim is as follows: If $H\mathbf{z} = \mathbf{0}$ where $\mathbf{z} \in P_{\mathbf{x}}$, then $HD\mathbf{f}(\mathbf{x})\mathbf{w} = \mathbf{0}$ where $\mathbf{w} \in M$ and $D\mathbf{f}(\mathbf{x})\mathbf{w} = \mathbf{z}$. Hence $\mathbf{w} = \mathbf{0}$ because $HD\mathbf{f}(\mathbf{x})$ is one to one, and so $\mathbf{z} = \mathbf{0}$ which shows that indeed H is one to one on $P_{\mathbf{x}}$.

Denote as $L_{\mathbf{x}}$ the inverse of H which is defined on $\mathbb{R}^m \times \mathbf{0}$, $L_{\mathbf{x}} : \mathbb{R}^m \times \mathbf{0} \rightarrow P_{\mathbf{x}}$. That $\mathbf{0}$ refers to the $N - m$ string of zeros in the definition given above for H .

Define $\mathbf{h} \equiv \mathbf{g}^{-1}$ and consider $\mathbf{f}_1 \equiv \mathbf{f} \circ \mathbf{h}$. It is desired to show that \mathbf{f}_1 depends only on x_1, \dots, x_m . Let D_1 refer to (x_1, \dots, x_m) and let D_2 refer to (x_{m+1}, \dots, x_n) . Then $\mathbf{f} = \mathbf{f}_1 \circ \mathbf{g}$ and so by the chain rule

$$D\mathbf{f}(\mathbf{x})(\mathbf{y}) = D\mathbf{f}_1(\mathbf{g}(\mathbf{x}))D\mathbf{g}(\mathbf{x})(\mathbf{y}) \quad (8.25)$$

Now as in 8.24, for $\mathbf{y} \in \mathbb{R}^n$,

$$\begin{aligned} D\mathbf{g}(\mathbf{x})(\mathbf{y}) &= HD\mathbf{f}(\mathbf{x})(\mathbf{y}) + G\mathbf{y} \\ &= (\phi_1(D\mathbf{f}(\mathbf{x})\mathbf{y}), \dots, \phi_m(D\mathbf{f}(\mathbf{x})\mathbf{y}), \psi_{m+1}(\mathbf{y}), \dots, \psi_n(\mathbf{y})) \end{aligned}$$

Recall that from the above definitions of H and G ,

$$\begin{aligned} G(\mathbf{y}) &\equiv (0, \dots, 0, \psi_{m+1}(\mathbf{y}), \dots, \psi_n(\mathbf{y})) \\ H(D\mathbf{f}(\mathbf{x})(\mathbf{y})) &= (\phi_1(D\mathbf{f}(\mathbf{x})\mathbf{y}), \dots, \phi_m(D\mathbf{f}(\mathbf{x})\mathbf{y}), 0, \dots, 0) \end{aligned}$$

Let $\pi_1 : \mathbb{R}^n \rightarrow \mathbb{R}^m$ denote the projection onto the first m positions and π_2 the projection onto the last $n - m$. Thus

$$\begin{aligned} \pi_1 D\mathbf{g}(\mathbf{x})(\mathbf{y}) &= (\phi_1(D\mathbf{f}(\mathbf{x})\mathbf{y}), \dots, \phi_m(D\mathbf{f}(\mathbf{x})\mathbf{y})) \\ \pi_2 D\mathbf{g}(\mathbf{x})(\mathbf{y}) &= (\psi_{m+1}(\mathbf{y}), \dots, \psi_n(\mathbf{y})) \end{aligned}$$

Now in general, for $\mathbf{z} \in \mathbb{R}^n$,

$$D\mathbf{f}_1(\mathbf{g}(\mathbf{x}))\mathbf{z} = D_1\mathbf{f}_1(\mathbf{g}(\mathbf{x}))\pi_1\mathbf{z} + D_2\mathbf{f}_1(\mathbf{g}(\mathbf{x}))\pi_2\mathbf{z}$$

Therefore, it follows that $D\mathbf{f}_1(\mathbf{g}(\mathbf{x}))D\mathbf{g}(\mathbf{x})(\mathbf{y})$ is given by

$$\begin{aligned} D\mathbf{f}(\mathbf{x})(\mathbf{y}) &= D\mathbf{f}_1(\mathbf{g}(\mathbf{x}))D\mathbf{g}(\mathbf{x})(\mathbf{y}) \\ &= D_1\mathbf{f}_1(\mathbf{g}(\mathbf{x}))\pi_1 D\mathbf{g}(\mathbf{x})(\mathbf{y}) + D_2\mathbf{f}_1(\mathbf{g}(\mathbf{x}))\pi_2 D\mathbf{g}(\mathbf{x})(\mathbf{y}) \end{aligned}$$

$$\begin{aligned}
Df(x)(y) &= Df_1(g(x))Dg(x)(y) = D_1f_1(g(x))\overbrace{\pi_1HDf(x)(y)}^{=\pi_1Dg(x)(y)} \\
&\quad + D_2f_1(g(x))\pi_2Gy
\end{aligned}$$

We need to verify the last term equals 0. Solving for this term,

$$D_2f_1(g(x))\pi_2Gy = Df(x)(y) - D_1f_1(g(x))\pi_1HDf(x)(y)$$

As just explained, $L_x \circ H$ is the identity on P_x , the image of $Df(x)$. Then

$$\begin{aligned}
D_2f_1(g(x))\pi_2Gy &= L_x \circ HDf(x)(y) - D_1f_1(g(x))\pi_1HDf(x)(y) \\
&= \left(L_x \circ \underline{HDf(x)} - D_1f_1(g(x))\pi_1\underline{HDf(x)} \right)(y)
\end{aligned}$$

Factoring out that underlined term,

$$D_2f_1(g(x))\pi_2Gy = [L_x - D_1f_1(g(x))\pi_1]HDf(x)(y)$$

Now $Df(x) : M \rightarrow P_x = Df(x)(\mathbb{R}^n)$ is onto. (This is based on the assumption that $Df(x)$ has rank m .) Thus it suffices to consider only $y \in M$ in the right side of the above. However, for such y , $\pi_2Gy = 0$ because to be in M , $\psi_k(y) = 0$ if $k \geq m+1$, and so the left side of the above equals 0 . Thus it appears this term on the left is 0 for any y chosen. How can this be so? It can only take place if $D_2f_1(g(x)) = 0$ for every $x \in V$. Thus, since g is onto, it can only take place if $D_2f_1(x) = 0$ for all $x \in U$. Therefore on U it must be the case that f_1 depends only on x_1, \dots, x_m as desired. ■

8.9 The Local Structure of C^1 Mappings

In linear algebra it is shown that every invertible matrix can be written as a product of elementary matrices, those matrices which are obtained from doing a row operation to the identity matrix. Two of the row operations produce a matrix which will change exactly one entry of a vector when it is multiplied by the elementary matrix. The other row operation involves switching two rows and this has the effect of switching two entries in a vector when multiplied on the left by the elementary matrix. Thus, in terms of the effect on a vector, the mapping determined by the given matrix can be considered as a composition of mappings which either flip two entries of the vector or change exactly one. A similar local result is available for nonlinear mappings. I found this interesting result in the advanced calculus book by Rudin.

Definition 8.9.1 Let U be an open set in \mathbb{R}^n and let $G : U \rightarrow \mathbb{R}^n$. Then G is called primitive if it is of the form

$$G(x) = \begin{pmatrix} x_1 & \cdots & \alpha(x) & \cdots & x_n \end{pmatrix}^T.$$

Thus, G is primitive if it only changes one of the variables. A function $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is called a flip if

$$F(x_1, \dots, x_k, \dots, x_l, \dots, x_n) = (x_1, \dots, x_l, \dots, x_k, \dots, x_n)^T.$$

Thus a function is a flip if it interchanges two coordinates. Also, for $m = 1, 2, \dots, n$, define

$$P_m(x) \equiv \begin{pmatrix} x_1 & x_2 & \cdots & x_m & 0 & \cdots & 0 \end{pmatrix}^T$$

It turns out that if $\mathbf{h}(\mathbf{0}) = \mathbf{0}$, $D\mathbf{h}(\mathbf{0})^{-1}$ exists, and \mathbf{h} is C^1 on U , then \mathbf{h} can be written as a composition of primitive functions and flips. This is a very interesting application of the inverse function theorem.

Theorem 8.9.2 *Let $\mathbf{h} : U \rightarrow \mathbb{R}^n$ be a C^1 function with $\mathbf{h}(\mathbf{0}) = \mathbf{0}$, $D\mathbf{h}(\mathbf{0})^{-1}$ exists. Then there is an open set $V \subseteq U$ containing $\mathbf{0}$, flips $\mathbf{F}_1, \dots, \mathbf{F}_{n-1}$, and primitive functions $\mathbf{G}_n, \mathbf{G}_{n-1}, \dots, \mathbf{G}_1$ such that for $\mathbf{x} \in V$,*

$$\mathbf{h}(\mathbf{x}) = \mathbf{F}_1 \circ \dots \circ \mathbf{F}_{n-1} \circ \mathbf{G}_n \circ \mathbf{G}_{n-1} \circ \dots \circ \mathbf{G}_1(\mathbf{x}).$$

The primitive function \mathbf{G}_j leaves x_i unchanged for $i \neq j$.

Proof: Let

$$\mathbf{h}_1(\mathbf{x}) \equiv \mathbf{h}(\mathbf{x}) = \begin{pmatrix} \alpha_1(\mathbf{x}) & \dots & \alpha_n(\mathbf{x}) \end{pmatrix}^T$$

$$D\mathbf{h}(\mathbf{0})\mathbf{e}_1 = \begin{pmatrix} \alpha_{1,1}(\mathbf{0}) & \dots & \alpha_{n,1}(\mathbf{0}) \end{pmatrix}^T$$

where $\alpha_{k,1}$ denotes $\frac{\partial \alpha_k}{\partial x_1}$. Since $D\mathbf{h}(\mathbf{0})$ is one to one, the right side of this expression cannot be zero. Hence there exists some k such that $\alpha_{k,1}(\mathbf{0}) \neq 0$. Now define

$$\mathbf{G}_1(\mathbf{x}) \equiv \begin{pmatrix} \alpha_k(\mathbf{x}) & x_2 & \dots & x_n \end{pmatrix}^T$$

Then the matrix of $D\mathbf{G}_1(\mathbf{0})$ is of the form

$$\begin{pmatrix} \alpha_{k,1}(\mathbf{0}) & \dots & \dots & \alpha_{k,n}(\mathbf{0}) \\ 0 & 1 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix}$$

and its determinant equals $\alpha_{k,1}(\mathbf{0}) \neq 0$. Therefore, by the inverse function theorem, there exists an open set U_1 , containing $\mathbf{0}$ and an open set V_2 containing $\mathbf{0}$ such that $\mathbf{G}_1(U_1) = V_2$ and \mathbf{G}_1 is one to one and onto, such that it and its inverse are both C^1 . Let \mathbf{F}_1 denote the flip which interchanges x_k with x_1 . Now define

$$\mathbf{h}_2(\mathbf{y}) \equiv \mathbf{F}_1 \circ \mathbf{h}_1 \circ \mathbf{G}_1^{-1}(\mathbf{y})$$

Thus

$$\begin{aligned} \mathbf{h}_2(\mathbf{G}_1(\mathbf{x})) &\equiv \mathbf{F}_1 \circ \mathbf{h}_1(\mathbf{x}) \\ &= \begin{pmatrix} \alpha_k(\mathbf{x}) & \dots & \alpha_1(\mathbf{x}) & \dots & \alpha_n(\mathbf{x}) \end{pmatrix}^T \end{aligned} \quad (8.26)$$

Therefore,

$$P_1 \mathbf{h}_2(\mathbf{G}_1(\mathbf{x})) = \begin{pmatrix} \alpha_k(\mathbf{x}) & 0 & \dots & 0 \end{pmatrix}^T.$$

Also

$$P_1(\mathbf{G}_1(\mathbf{x})) = \begin{pmatrix} \alpha_k(\mathbf{x}) & 0 & \dots & 0 \end{pmatrix}^T$$

so $P_1 \mathbf{h}_2(\mathbf{y}) = P_1(\mathbf{y})$ for all $\mathbf{y} \in V_2$. Also, $\mathbf{h}_2(\mathbf{0}) = \mathbf{0}$ and $D\mathbf{h}_2(\mathbf{0})^{-1}$ exists because of the definition of \mathbf{h}_2 above and the chain rule. Since $\mathbf{F}_1^2 = I$, the identity map, it follows from (8.26) that

$$\mathbf{h}(\mathbf{x}) = \mathbf{h}_1(\mathbf{x}) = \mathbf{F}_1 \circ \mathbf{h}_2 \circ \mathbf{G}_1(\mathbf{x}). \quad (8.27)$$

Note that on an open set $V_2 \equiv G_1(U_1)$ containing the origin, h_2 leaves the first entry unchanged. This is what $P_1 h_2(G_1(x)) = P_1(G_1(x))$ says. In contrast, $h_1 = h$ left possibly no entries unchanged.

Suppose then, that for $m \geq 2$, h_m leaves the first $m - 1$ entries unchanged,

$$P_{m-1} h_m(x) = P_{m-1}(x) \quad (8.28)$$

for all $x \in U_m$, an open subset of U containing 0 , and $h_m(0) = 0$, $Dh_m(0)^{-1}$ exists. From (8.28), $h_m(x)$ must be of the form

$$h_m(x) = \begin{pmatrix} x_1 & \cdots & x_{m-1} & \alpha_1(x) & \cdots & \alpha_n(x) \end{pmatrix}^T$$

where these α_k are different than the ones used earlier. Then

$$Dh_m(0) e_m = \begin{pmatrix} 0 & \cdots & 0 & \alpha_{1,m}(0) & \cdots & \alpha_{n,m}(0) \end{pmatrix}^T \neq 0$$

because $Dh_m(0)^{-1}$ exists. Therefore, there exists a $k \geq m$ such that $\alpha_{k,m}(0) \neq 0$, not the same k as before. Define

$$G_m(x) \equiv \begin{pmatrix} x_1 & \cdots & x_{m-1} & \alpha_k(x) & \cdots & x_n \end{pmatrix}^T \quad (8.29)$$

so a change in G_m occurs only in the m^{th} slot. Then $G_m(0) = 0$ and $DG_m(0)^{-1}$ exists similar to the above. In fact

$$\det(DG_m(0)) = \alpha_{k,m}(0).$$

Therefore, by the inverse function theorem, there exists an open set V_{m+1} containing 0 such that $V_{m+1} = G_m(U_m)$ with G_m and its inverse being one to one, continuous and onto. Let F_m be the flip which flips x_m and x_k . Then define h_{m+1} on V_{m+1} by

$$h_{m+1}(y) = F_m \circ h_m \circ G_m^{-1}(y).$$

Thus for $x \in U_m$,

$$h_{m+1}(G_m(x)) = (F_m \circ h_m)(x). \quad (8.30)$$

and consequently, since $F_m^2 = I$,

$$F_m \circ h_{m+1} \circ G_m(x) = h_m(x) \quad (8.31)$$

It follows

$$\begin{aligned} P_m h_{m+1}(G_m(x)) &= P_m(F_m \circ h_m)(x) \\ &= \begin{pmatrix} x_1 & \cdots & x_{m-1} & \alpha_k(x) & 0 & \cdots & 0 \end{pmatrix}^T \end{aligned}$$

and

$$P_m(G_m(x)) = \begin{pmatrix} x_1 & \cdots & x_{m-1} & \alpha_k(x) & 0 & \cdots & 0 \end{pmatrix}^T.$$

Therefore, for $y \in V_{m+1}$,

$$P_m h_{m+1}(y) = P_m(y).$$

As before, $h_{m+1}(\mathbf{0}) = \mathbf{0}$ and $Dh_{m+1}(\mathbf{0})^{-1}$ exists. Therefore, we can apply (8.31) repeatedly, obtaining the following:

$$\begin{aligned} h(x) &= F_1 \circ h_2 \circ G_1(x) \\ &= F_1 \circ F_2 \circ h_3 \circ G_2 \circ G_1(x) \\ &\vdots \\ &= F_1 \circ \cdots \circ F_{n-1} \circ h_n \circ G_{n-1} \circ \cdots \circ G_1(x) \end{aligned}$$

where h_n fixes the first $n-1$ entries,

$$P_{n-1}h_n(x) = P_{n-1}(x) = \begin{pmatrix} x_1 & \cdots & x_{n-1} & 0 \end{pmatrix}^T,$$

and so $h_n(x)$ is a primitive mapping of the form

$$h_n(x) = \begin{pmatrix} x_1 & \cdots & x_{n-1} & \alpha(x) \end{pmatrix}^T.$$

Therefore, define the primitive function $G_n(x)$ to equal $h_n(x)$. ■

8.10 Invariance of Domain

As an application of the inverse function theorem is a simple proof of the important invariance of domain theorem which says that continuous and one to one functions defined on an open set in \mathbb{R}^n with values in \mathbb{R}^n take open sets to open sets. You know that this is true for functions of one variable because a one to one continuous function must be either strictly increasing or strictly decreasing. This will be used when considering orientations of curves later. However, the n dimensional version isn't at all obvious but is just as important if you want to consider manifolds with boundary for example. The need for this theorem occurs in many other places as well in addition to being extremely interesting for its own sake. The inverse function theorem gives conditions under which a differentiable function maps open sets to open sets. The following lemma, depending on the Brouwer fixed point theorem is the thing which will allow this to be extended to continuous one to one functions. It says roughly that if a continuous function does not move points near p very far, then the image of a ball centered at p contains an open set.

Lemma 8.10.1 *Let f be continuous and map $\overline{B(p, r)} \subseteq \mathbb{R}^n$ to \mathbb{R}^n . Suppose that for all $x \in \overline{B(p, r)}$, $|f(x) - x| < \varepsilon r$. Then it follows that $f(\overline{B(p, r)}) \supseteq B(p, (1 - \varepsilon)r)$*

Proof: This is from the Brouwer fixed point theorem, Corollary 6.3.2. Consider for $y \in B(p, (1 - \varepsilon)r)$,

$$h(x) \equiv x - f(x) + y$$

Then h is continuous and for $x \in \overline{B(p, r)}$,

$$|h(x) - p| = |x - f(x) + y - p| < \varepsilon r + |y - p| < \varepsilon r + (1 - \varepsilon)r = r$$

Hence $h : \overline{B(p, r)} \rightarrow \overline{B(p, r)}$ and so it has a fixed point x by Corollary 6.3.2 or Theorem 11.6.8. Thus

$$x - f(x) + y = x$$

so $f(x) = y$. ■

The notation $\|f\|_K$ means $\sup_{x \in K} |f(x)|$. If you have a continuous function h defined on a compact set K , then the Stone Weierstrass theorem implies you can uniformly approximate it with a polynomial g . That is $\|h - g\|_K$ is small. The following lemma says that you can also have $g(z) = h(z)$ and $Dg(z)^{-1}$ exists so that near z , the function g will map open sets to open sets as claimed by the inverse function theorem. First is a little observation about approximating.

Lemma 8.10.2 *Suppose $\det(A) = 0$. Then for all sufficiently small nonzero ε ,*

$$\det(A + \varepsilon I) \neq 0$$

Proof: First suppose A is a $p \times p$ matrix. Suppose also that $\det(A) = 0$. Thus, the constant term of $\det(\lambda I - A)$ is 0. Consider $\varepsilon I + A \equiv A_\varepsilon$ for small real ε . The characteristic polynomial of A_ε is

$$\det(\lambda I - A_\varepsilon) = \det((\lambda - \varepsilon)I - A)$$

This is of the form

$$(\lambda - \varepsilon)^p + a_{p-1}(\lambda - \varepsilon)^{p-1} + \cdots + (\lambda - \varepsilon)^m a_m$$

where the a_j are the coefficients in the characteristic equation for A and m is the largest such that $a_m \neq 0$. The constant term of this characteristic polynomial for A_ε must be nonzero for all ε small enough because it is of the form

$$(-1)^m \varepsilon^m a_m + (\text{higher order terms in } \varepsilon)$$

which shows that $\varepsilon I + A$ is invertible for all ε small enough but nonzero. ■

Lemma 8.10.3 *Let K be a compact set in \mathbb{R}^n and let $h : K \rightarrow \mathbb{R}^n$ be continuous, $z \in K$ is fixed. Let $\delta > 0$. Then there exists a polynomial g (each component a polynomial) such that*

$$\|g - h\|_K < \delta, \quad g(z) = h(z), \quad Dg(z)^{-1} \text{ exists}$$

Proof: By the Weierstrass approximation theorem, Corollary 5.8.8, or Theorem 5.10.5, there exists a polynomial \hat{g} such that

$$\|\hat{g} - h\|_K < \frac{\delta}{3}$$

Then define for $y \in K$

$$g(y) \equiv \hat{g}(y) + h(z) - \hat{g}(z)$$

Then

$$g(z) = \hat{g}(z) + h(z) - \hat{g}(z) = h(z)$$

Also

$$\begin{aligned} |g(y) - h(y)| &\leq |(\hat{g}(y) + h(z) - \hat{g}(z)) - h(y)| \\ &\leq |\hat{g}(y) - h(y)| + |h(z) - \hat{g}(z)| < \frac{2\delta}{3} \end{aligned}$$

and so since y was arbitrary,

$$\|g - h\|_K \leq \frac{2\delta}{3} < \delta$$

If $Dg(z)^{-1}$ exists, then this is what is wanted. If not, use Lemma 8.10.2 and note that for all η small enough, you could replace g with $y \rightarrow g(y) + \eta(y - z)$ and it will still be the case that $\|g - h\|_K < \delta$ along with $g(z) = h(z)$ but now $Dg(z)^{-1}$ exists. Simply use the modified g . ■

The main result is essentially the following lemma which combines the conclusions of the above.

Lemma 8.10.4 *Let $f : \overline{B(p, r)} \rightarrow \mathbb{R}^n$ where the ball is also in \mathbb{R}^n . Let f be one to one, f continuous. Then there exists $\delta > 0$ such that*

$$f(\overline{B(p, r)}) \supseteq B(f(p), \delta).$$

In other words, $f(p)$ is an interior point of $f(\overline{B(p, r)})$.

Proof: Since $f(\overline{B(p, r)})$ is compact, it follows that $f^{-1} : f(\overline{B(p, r)}) \rightarrow \overline{B(p, r)}$ is continuous. By Lemma 8.10.3, there exists a polynomial $g : f(\overline{B(p, r)}) \rightarrow \mathbb{R}^n$ such that

$$\begin{aligned} \|g - f^{-1}\|_{f(\overline{B(p, r)})} &< \varepsilon r, \varepsilon < 1, \quad Dg(f(p))^{-1} \\ \text{exists, and } g(f(p)) &= f^{-1}(f(p)) = p \end{aligned}$$

From the first inequality in the above,

$$|g(f(x)) - x| = |g(f(x)) - f^{-1}(f(x))| \leq \|g - f^{-1}\|_{f(\overline{B(p, r)})} < \varepsilon r$$

By Lemma 8.10.1,

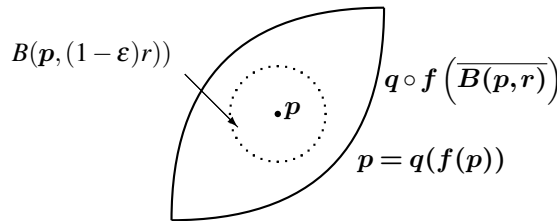
$$g \circ f(\overline{B(p, r)}) \supseteq B(p, (1 - \varepsilon)r) = B(g(f(p)), (1 - \varepsilon)r)$$

Since $Dg(f(p))^{-1}$ exists, it follows from the inverse function theorem that g^{-1} also exists and that g, g^{-1} are open maps on small open sets containing $f(p)$ and p respectively. Thus there exists $\eta < (1 - \varepsilon)r$ such that g^{-1} is an open map on $B(p, \eta) \subseteq B(p, (1 - \varepsilon)r)$. Thus

$$g \circ f(\overline{B(p, r)}) \supseteq B(p, (1 - \varepsilon)r) \supseteq B(p, \eta)$$

So do g^{-1} to both ends. Then you have $g^{-1}(p) = f(p)$ is in the open set $g^{-1}(B(p, \eta))$. Thus

$$f(\overline{B(p, r)}) \supseteq g^{-1}(B(p, \eta)) \supseteq B(g^{-1}(p), \delta) = B(f(p), \delta) \quad \blacksquare$$



With this lemma, the invariance of domain theorem comes right away. This remarkable theorem states that if $f : U \rightarrow \mathbb{R}^n$ for U an open set in \mathbb{R}^n and if f is one to one and continuous, then $f(U)$ is also an open set in \mathbb{R}^n .

Theorem 8.10.5 *Let U be an open set in \mathbb{R}^n and let $f : U \rightarrow \mathbb{R}^n$ be one to one and continuous. Then $f(U)$ is also an open subset in \mathbb{R}^n .*

Proof: It suffices to show that if $p \in U$ then $f(p)$ is an interior point of $f(U)$. Let $\overline{B(p, r)} \subseteq U$. By Lemma 8.10.4, $f(U) \supseteq f(\overline{B(p, r)}) \supseteq B(f(p), \delta)$ so $f(p)$ is indeed an interior point of $f(U)$. ■

The inverse mapping theorem assumed quite a bit about the mapping. In particular it assumed that the mapping had a continuous derivative. The following version of the inverse function theorem seems very interesting because it only needs an invertible derivative at a point.

Corollary 8.10.6 *Let U be an open set in \mathbb{R}^p and let $f : U \rightarrow \mathbb{R}^p$ be one to one and continuous. Then, f^{-1} is also continuous on the open set $f(U)$. If f is differentiable at $x_1 \in U$ and if $Df(x_1)^{-1}$ exists for $x_1 \in U$, then it follows that $Df^{-1}(f(x_1)) = Df(x_1)^{-1}$.*

Proof: $|\cdot|$ will be a norm on \mathbb{R}^p , whichever is desired. If you like, let it be the Euclidean norm. $\|\cdot\|$ will be the operator norm. The first part of the conclusion of this corollary is from invariance of domain.

From the assumption that $Df(x_1)$ and $Df(x_1)^{-1}$ exists,

$$y - f(x_1) = f(f^{-1}(y)) - f(x_1) = Df(x_1)(f^{-1}(y) - x_1) + o(f^{-1}(y) - x_1)$$

Since $Df(x_1)^{-1}$ exists,

$$Df(x_1)^{-1}(y - f(x_1)) = f^{-1}(y) - x_1 + o(f^{-1}(y) - x_1)$$

by continuity, if $|y - f(x_1)|$ is small enough, then $|f^{-1}(y) - x_1|$ is small enough that in the above,

$$|o(f^{-1}(y) - x_1)| < \frac{1}{2}|f^{-1}(y) - x_1|$$

Hence, if $|y - f(x_1)|$ is sufficiently small, then from the triangle inequality of the form $|p - q| \geq ||p| - |q||$,

$$\begin{aligned} \left\| Df(x_1)^{-1} \right\| |(y - f(x_1))| &\geq \left| Df(x_1)^{-1}(y - f(x_1)) \right| \\ &\geq |f^{-1}(y) - x_1| - \frac{1}{2}|f^{-1}(y) - x_1| = \frac{1}{2}|f^{-1}(y) - x_1| \\ |y - f(x_1)| &\geq \left\| Df(x_1)^{-1} \right\|^{-1} \frac{1}{2}|f^{-1}(y) - x_1| \end{aligned}$$

It follows that for $|y - f(x_1)|$ small enough,

$$\left| \frac{o(f^{-1}(y) - x_1)}{y - f(x_1)} \right| \leq \left| \frac{o(f^{-1}(y) - x_1)}{f^{-1}(y) - x_1} \right| \frac{2}{\left\| Df(x_1)^{-1} \right\|^{-1}}$$

Then, using continuity of the inverse function again, it follows that if $|y - f(x_1)|$ is possibly still smaller, then $f^{-1}(y) - x_1$ is sufficiently small that the right side of the

above inequality is no larger than ε . Since ε is arbitrary, it follows $\mathbf{o}(\mathbf{f}^{-1}(\mathbf{y}) - \mathbf{x}_1) = \mathbf{o}(\mathbf{y} - \mathbf{f}(\mathbf{x}_1))$. Now from differentiability of \mathbf{f} at \mathbf{x}_1 ,

$$\begin{aligned} \mathbf{y} - \mathbf{f}(\mathbf{x}_1) &= \mathbf{f}(\mathbf{f}^{-1}(\mathbf{y})) - \mathbf{f}(\mathbf{x}_1) = D\mathbf{f}(\mathbf{x}_1)(\mathbf{f}^{-1}(\mathbf{y}) - \mathbf{x}_1) + \mathbf{o}(\mathbf{f}^{-1}(\mathbf{y}) - \mathbf{x}_1) \\ &= D\mathbf{f}(\mathbf{x}_1)(\mathbf{f}^{-1}(\mathbf{y}) - \mathbf{x}_1) + \mathbf{o}(\mathbf{y} - \mathbf{f}(\mathbf{x}_1)) \\ &= D\mathbf{f}(\mathbf{x}_1)(\mathbf{f}^{-1}(\mathbf{y}) - \mathbf{f}^{-1}(\mathbf{f}(\mathbf{x}_1))) + \mathbf{o}(\mathbf{y} - \mathbf{f}(\mathbf{x}_1)) \end{aligned}$$

Therefore, solving for $\mathbf{f}^{-1}(\mathbf{y}) - \mathbf{f}^{-1}(\mathbf{f}(\mathbf{x}_1))$,

$$\mathbf{f}^{-1}(\mathbf{y}) - \mathbf{f}^{-1}(\mathbf{f}(\mathbf{x}_1)) = D\mathbf{f}(\mathbf{x}_1)^{-1}(\mathbf{y} - \mathbf{f}(\mathbf{x}_1)) + \mathbf{o}(\mathbf{y} - \mathbf{f}(\mathbf{x}_1))$$

From the definition of the derivative, this shows that $D\mathbf{f}^{-1}(\mathbf{f}(\mathbf{x}_1)) = D\mathbf{f}(\mathbf{x}_1)^{-1}$. ■

8.11 Exercises

1. This problem was suggested to me by Matt Heiner. Earlier there was a problem in which two surfaces intersected at a point and this implied that in fact, they intersected in a smooth curve. Now suppose you have two spheres $x^2 + y^2 + z^2 = 1$ and $(x-2)^2 + y^2 + z^2 = 1$. These intersect at the single point $(1, 0, 0)$. Why does the implicit function theorem not imply that these surfaces intersect in a curve?
2. Maximize $2x + y$ subject to the condition that $\frac{x^2}{4} + \frac{y^2}{9} \leq 1$. **Hint:** You need to consider interior points and also the method of Lagrange multipliers for the points on the boundary of this ellipse.
3. Maximize $x + y$ subject to the condition that $x^2 + \frac{y^2}{9} + z^2 \leq 1$.
4. Find the points on $y^2x = 16$ which are closest to $(0, 0)$.
5. Use Lagrange multipliers to “solve” the following maximization problem. Maximize xy^2z^3 subject to the constraint $x + y + z = 12$. Show that the Lagrange multiplier method works very well but gives an answer which is neither a maximum nor a minimum. **Hint:** Show there is no maximum by considering $y = 12 - 5x, z = 4x$ and then letting x be large.
6. Let $f(x, y, z) = x^2 - 2yx + 2z^2 - 4z + 2$. Identify all the points where $Df = 0$. Then determine whether they are local minima local maxima or saddle points.
7. Let $f(x, y) = x^4 - 2x^2 + 2y^2 + 1$. Identify all the points where $Df = 0$. Then determine whether they are local minima local maxima or saddle points.
8. Let $f(x, y, z) = -x^4 + 2x^2 - y^2 - 2z^2 - 1$. Identify all the points where $Df = 0$. Then determine whether they are local minima local maxima or saddle points.
9. Let $f : V \rightarrow \mathbb{R}$ where V is a finite dimensional normed vector space. Suppose f is convex which means $f(t\mathbf{x} + (1-t)\mathbf{y}) \leq tf(\mathbf{x}) + (1-t)f(\mathbf{y})$ whenever $t \in [0, 1]$. Suppose also that f is differentiable. Show then that for every $\mathbf{x}, \mathbf{y} \in V$,

$$(Df(\mathbf{x}) - Df(\mathbf{y}))(\mathbf{x} - \mathbf{y}) \geq 0.$$

Thus convex functions have monotone derivatives.

10. Suppose B is an open ball in X and $\mathbf{f} : B \rightarrow Y$ is differentiable. Suppose also there exists $L \in \mathcal{L}(X, Y)$ such that $\|D\mathbf{f}(\mathbf{x}) - L\| < k$ for all $\mathbf{x} \in B$. Show that if $\mathbf{x}_1, \mathbf{x}_2 \in B$,

$$\|\mathbf{f}(\mathbf{x}_1) - \mathbf{f}(\mathbf{x}_2) - L(\mathbf{x}_1 - \mathbf{x}_2)\| \leq k \|\mathbf{x}_1 - \mathbf{x}_2\|.$$

Hint: Consider $T\mathbf{x} = \mathbf{f}(\mathbf{x}) - L\mathbf{x}$ and argue $\|DT(\mathbf{x})\| < k$.

11. Let $\mathbf{f} : U \subseteq X \rightarrow Y$, $D\mathbf{f}(\mathbf{x})$ exists for all $\mathbf{x} \in U$, $B(\mathbf{x}_0, \delta) \subseteq U$, and there exists $L \in \mathcal{L}(X, Y)$, such that $L^{-1} \in \mathcal{L}(Y, X)$, and for all $\mathbf{x} \in B(\mathbf{x}_0, \delta)$

$$\|D\mathbf{f}(\mathbf{x}) - L\| < \frac{r}{\|L^{-1}\|}, \quad r < 1.$$

Show that there exists $\varepsilon > 0$ and an open subset of $B(\mathbf{x}_0, \delta)$ called V , such that $\mathbf{f} : V \rightarrow B(\mathbf{f}(\mathbf{x}_0), \varepsilon)$ is one to one and onto. Also $D\mathbf{f}^{-1}(\mathbf{y})$ exists for each $\mathbf{y} \in B(\mathbf{f}(\mathbf{x}_0), \varepsilon)$ and is given by the formula $D\mathbf{f}^{-1}(\mathbf{y}) = [D\mathbf{f}(\mathbf{f}^{-1}(\mathbf{y}))]^{-1}$. **Hint:** Let

$$T_{\mathbf{y}}(\mathbf{x}) \equiv T(\mathbf{x}, \mathbf{y}) \equiv \mathbf{x} - L^{-1}(\mathbf{f}(\mathbf{x}) - \mathbf{y})$$

for $\|\mathbf{y} - \mathbf{f}(\mathbf{x}_0)\| < \frac{(1-r)\delta}{2\|L^{-1}\|}$, consider $\{T_{\mathbf{y}}^n(\mathbf{x}_0)\}$. This is a version of the inverse function theorem for \mathbf{f} only differentiable, not C^1 .

12. If \mathbf{f} is one to one and C^1 , and $D\mathbf{f}(\mathbf{x}_0)$ is invertible, then locally the function \mathbf{f} is one to one. Explain why this is, maybe using the above problem. However, this is a strictly local result! Let $\mathbf{f} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ be given by

$$\mathbf{f}(x, y) = \begin{pmatrix} e^x \cos y \\ e^x \sin y \end{pmatrix}$$

This clearly is not one to one because if you replace y with $y + 2\pi$, you get the same value. Now verify that $D\mathbf{f}(x, y)^{-1}$ exists for all (x, y) .

13. Show every polynomial, $\sum_{|\alpha| \leq k} d_{\alpha} \mathbf{x}^{\alpha}$ is C^k for every k . Show that if f is defined and continuous on a compact set K , then there is an infinitely differentiable function which is uniformly close to f on K .
14. Suppose $U \subseteq \mathbb{R}^2$ is an open set and $\mathbf{f} : U \rightarrow \mathbb{R}^3$ is C^1 . Suppose $D\mathbf{f}(s_0, t_0)$ has rank two and

$$\mathbf{f}(s_0, t_0) = \begin{pmatrix} x_0 & y_0 & z_0 \end{pmatrix}^T.$$

Show that for (s, t) near (s_0, t_0) , the points $\mathbf{f}(s, t)$ may be realized in one of the following forms.

$$\{(x, y, \phi(x, y)) : (x, y) \text{ near } (x_0, y_0)\},$$

$$\{(\phi(y, z), y, z) : (y, z) \text{ near } (y_0, z_0)\},$$

or

$$\{(x, \phi(x, z), z) : (x, z) \text{ near } (x_0, z_0)\}.$$

This shows that parametrically defined surfaces can be obtained locally in a particularly simple form.

15. Minimize $\sum_{j=1}^n x_j$ subject to the constraint $\sum_{j=1}^n x_j^2 = a^2$. Your answer should be some function of a which you may assume is a positive number.

16. A curve is formed from the intersection of the plane, $2x + 3y + z = 3$ and the cylinder $x^2 + y^2 = 4$. Find the point on this curve which is closest to $(0, 0, 0)$.
17. A curve is formed from the intersection of the plane, $2x + 3y + z = 3$ and the sphere $x^2 + y^2 + z^2 = 16$. Find the point on this curve which is closest to $(0, 0, 0)$.
18. Let $A = (A_{ij})$ be an $n \times n$ matrix which is symmetric. Thus $A_{ij} = A_{ji}$ and recall $(A\mathbf{x})_i = A_{ij}x_j$ where you sum over the repeated index. Show $\frac{\partial}{\partial x_i} (A_{ij}x_jx_i) = 2A_{ij}x_j$. Show that when you use the method of Lagrange multipliers to maximize the function, $A_{ij}x_jx_i$ subject to the constraint, $\sum_{j=1}^n x_j^2 = 1$, the value of λ which corresponds to the maximum value of this functions is such that $A_{ij}x_j = \lambda x_i$. Thus $A\mathbf{x} = \lambda\mathbf{x}$. Thus λ is an eigenvalue of the matrix A .
19. Let x_1, \dots, x_5 be 5 positive numbers. Maximize their product subject to the constraint that

$$x_1 + 2x_2 + 3x_3 + 4x_4 + 5x_5 = 300.$$

20. Let $f(x_1, \dots, x_n) = x_1^n x_2^{n-1} \cdots x_n^1$. Then f achieves a maximum on the set,

$$S \equiv \left\{ \mathbf{x} \in \mathbb{R}^n : \sum_{i=1}^n ix_i = 1 \text{ and each } x_i \geq 0 \right\}.$$

If $\mathbf{x} \in S$ is the point where this maximum is achieved, find x_1/x_n .

21. Maximize $\prod_{i=1}^n x_i^2$ subject to the constraint, $\sum_{i=1}^n x_i^2 = r^2$. Show the maximum is $(r^2/n)^n$. Now show from this that $(\prod_{i=1}^n x_i^2)^{1/n} \leq \frac{1}{n} \sum_{i=1}^n x_i^2$ and finally, conclude that if each number $x_i \geq 0$, then

$$\left(\prod_{i=1}^n x_i \right)^{1/n} \leq \frac{1}{n} \sum_{i=1}^n x_i$$

and there exist values of the x_i for which equality holds. This says the “geometric mean” is always smaller than the arithmetic mean.

22. Show that there exists a smooth solution $y = y(x)$ to the equation

$$xe^y + ye^x = 0$$

which is valid for x, y both near 0. Find $y'(x)$ at a point (x, y) near $(0, 0)$. Then find $y''(x)$ for such (x, y) . Can you find an explicit formula for $y(x)$?

23. The next few problems involve invariance of domain. Suppose U is a nonempty open set in \mathbb{R}^n , $f: U \rightarrow \mathbb{R}^n$ is continuous, and suppose that for each $\mathbf{x} \in U$, there is a ball $B_{\mathbf{x}}$ containing \mathbf{x} such that f is one to one on $B_{\mathbf{x}}$. That is, f is locally one to one. Show that $f(U)$ is open.
24. \uparrow In the situation of the above problem, suppose $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ is locally one to one. Also suppose that $\lim_{|\mathbf{x}| \rightarrow \infty} |f(\mathbf{x})| = \infty$. Show it follows that $f(\mathbb{R}^n) = \mathbb{R}^n$. That is, f is onto. Show that this would not be true if f is only defined on a proper open set. Also show that this would not be true if the condition $\lim_{|\mathbf{x}| \rightarrow \infty} |f(\mathbf{x})| = \infty$ does not hold. **Hint:** You might show that $f(\mathbb{R}^n)$ is both open and closed and then use connectedness. To get an example in the second case, you might think of e^{x+iy} . It does not include $0 + i0$. Why not?

25. \uparrow Show that if $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is C^1 and if $D\mathbf{f}(\mathbf{x})$ exists and is invertible for all $\mathbf{x} \in \mathbb{R}^n$, then \mathbf{f} is locally one to one. Thus, from the above problem, if $\lim_{|\mathbf{x}| \rightarrow \infty} |\mathbf{f}(\mathbf{x})| = \infty$, then \mathbf{f} is also onto. Now consider $\mathbf{f} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ given by

$$\mathbf{f}(x, y) = \begin{pmatrix} e^x \cos y \\ e^x \sin y \end{pmatrix}$$

Show that this does not map onto \mathbb{R}^2 . In fact, it fails to hit $(0, 0)$, but $D\mathbf{f}(x, y)$ is invertible for all (x, y) . Show why it fails to satisfy the limit condition.

26. You know from linear algebra that there is no onto linear mapping $A : \mathbb{R}^m \rightarrow \mathbb{R}^p$ for $p > m$. Show that there is no locally one to one continuous mapping which will map \mathbb{R}^m onto \mathbb{R}^p .
27. In Example 8.1.9 on Page 210, could you replace y with $\mathbf{y} \in \mathbb{R}^m$ and obtain a modified version of this example?

Part II

Integration

Chapter 9

Measures and Measurable Functions

The Lebesgue integral is much better than the Riemann integral. This has been known for over 100 years. It is **much easier** to generalize to many dimensions and it is much easier to use in applications. It is also this integral which is most important in probability. However, this integral is more abstract. This chapter will develop the abstract machinery for this integral.

The next definition describes what is meant by a σ algebra. This is the fundamental object which is studied in probability theory. The events come from a σ algebra of sets. Recall that $\mathcal{P}(\Omega)$ is the set of all subsets of the given set Ω . It may also be denoted by 2^Ω but I won't refer to it this way.

Definition 9.0.1 $\mathcal{F} \subseteq \mathcal{P}(\Omega)$, the set of all subsets of Ω , is called a σ algebra if it contains \emptyset, Ω , and is closed with respect to countable unions and complements. That is, if $\{A_n\}_{n=1}^\infty$ is countable and each $A_n \in \mathcal{F}$, then $\bigcup_{n=1}^\infty A_n \in \mathcal{F}$ also and if $A \in \mathcal{F}$, then $\Omega \setminus A \in \mathcal{F}$. It is clear that any intersection of σ algebras is a σ algebra. If $\mathcal{K} \subseteq \mathcal{P}(\Omega)$, $\sigma(\mathcal{K})$ is the smallest σ algebra which contains \mathcal{K} . In fact, the intersection of all σ algebras containing \mathcal{K} is obviously a σ algebra so this intersection is $\sigma(\mathcal{K})$.

If \mathcal{F} is a σ algebra, then it is also closed with respect to countable intersections. Here is why. Let $\{F_k\}_{k=1}^\infty \subseteq \mathcal{F}$. Then $(\bigcap_k F_k)^C = \bigcup_k F_k^C \in \mathcal{F}$ and so $\bigcap_k F_k = \left(\left(\bigcap_k F_k\right)^C\right)^C = \left(\bigcup_k F_k^C\right)^C \in \mathcal{F}$.

Example 9.0.2 You could consider \mathbb{N} and for your σ algebra, you could have $\mathcal{P}(\mathbb{N})$. This satisfies all the necessary requirements. Note that in fact, $\mathcal{P}(S)$ works for any S . However, useful examples are not typically the set of all subsets.

9.1 Simple Functions and Measurable Functions

A σ algebra is a collection of subsets of a set Ω which includes \emptyset, Ω , and is closed with respect to countable unions and complements.

Definition 9.1.1 A measurable space, denoted as (Ω, \mathcal{F}) , is one for which \mathcal{F} is a σ algebra contained in $\mathcal{P}(\Omega)$. Let $f : \Omega \rightarrow X$ where X is a metric space. Then f is said to be measurable means $f^{-1}(U) \in \mathcal{F}$ whenever U is open.

It is important to have a theorem about pointwise limits of measurable functions. The following is a fairly general such theorem which holds in the situations to be considered in this book. First recall $\text{dist}(x, S)$ in Lemma 3.12.1 which implies that $x \rightarrow \text{dist}(x, S)$ is continuous.

Theorem 9.1.2 Let $\{f_n\}$ be a sequence of measurable functions mapping Ω to the metric space (X, d) where (Ω, \mathcal{F}) is a measurable space. Suppose the pointwise limit $f(\omega) = \lim_{n \rightarrow \infty} f_n(\omega)$ for all ω . Then f is also a measurable function.

Proof: It is required to show $f^{-1}(U)$ is measurable for all U open. Let

$$V_m \equiv \left\{ x \in U : \text{dist}(x, U^C) > \frac{1}{m} \right\}.$$

Thus, since dist is continuous, (Lemma 3.12.1), $V_m \subseteq \{x \in U : \text{dist}(x, U^C) \geq \frac{1}{m}\}$, $V_m \subseteq \overline{V_m} \subseteq V_{m+1}$, and $\cup_m V_m = U$. Then since V_m is open, $f^{-1}(V_m) = \cup_{n=1}^{\infty} \cap_{k=n}^{\infty} f_k^{-1}(V_m)$ and so

$$\begin{aligned} f^{-1}(U) &= \cup_{m=1}^{\infty} f^{-1}(V_m) = \cup_{m=1}^{\infty} \cup_{n=1}^{\infty} \cap_{k=n}^{\infty} f_k^{-1}(V_m) \\ &\subseteq \cup_{m=1}^{\infty} f^{-1}(\overline{V_m}) = f^{-1}(U) \end{aligned}$$

which shows $f^{-1}(U)$ is measurable. ■

Important examples of a metric spaces are $\mathbb{R}, \mathbb{C}, \mathbb{F}^n$, where \mathbb{F} is either \mathbb{R} or \mathbb{C} . However, it is also very convenient to consider the metric space $(-\infty, \infty]$, the real line with ∞ tacked on at the end. This can be considered as a metric space in a very simple way.

$$\rho(x, y) = |\arctan(x) - \arctan(y)|$$

with the understanding that $\arctan(\infty) \equiv \pi/2$. It is easy to show that this metric restricted to \mathbb{R} gives the same open sets on \mathbb{R} as the usual metric given by $d(x, y) = |x - y|$ but in addition, allows the inclusion of that ideal point out at the end of the real line denoted as ∞ . This is considered mainly because it makes the development of the theory easier. The open sets in $(-\infty, \infty]$ are described in the following lemma.

Lemma 9.1.3 *The open balls in $(-\infty, \infty]$ consist of sets of the form (a, b) for a, b real numbers and $(a, \infty]$. This is a separable metric space.*

Proof: If the center of the ball is a real number, then the ball will result in an interval (a, b) where a, b are real numbers. If the center of the ball is ∞ , then the ball results in something of the form $(a, \infty]$. It is obvious that this is a separable metric space with the countable dense set being \mathbb{Q} since every ball contains a rational number. ■

If you kept both $-\infty$ and ∞ with the obvious generalization that $\arctan(-\infty) \equiv -\frac{\pi}{2}$, then the resulting metric space would be a complete separable metric space. However, it is not convenient to include $-\infty$, so this won't be done. The reason is that it will be desired to make sense of things like $f + g$.

Then for functions which have values in $(-\infty, \infty]$ we have the following extremely useful description of what it means for a function to be measurable.

Lemma 9.1.4 *Let $f : \Omega \rightarrow (-\infty, \infty]$ where \mathcal{F} is a σ algebra of subsets of Ω . Here $(-\infty, \infty]$ is the metric space just described with the metric given by*

$$\rho(x, y) = |\arctan(x) - \arctan(y)|.$$

Then the following are equivalent.

$$\begin{aligned} f^{-1}((d, \infty]) &\in \mathcal{F}, \text{ for all finite } d, \\ f^{-1}((-\infty, d)) &\in \mathcal{F}, \text{ for all finite } d, \\ f^{-1}([d, \infty]) &\in \mathcal{F}, \text{ for all finite } d, \\ f^{-1}((-\infty, d]) &\in \mathcal{F}, \text{ for all finite } d, \\ f^{-1}((a, b)) &\in \mathcal{F} \text{ for all } a < b, -\infty < a < b < \infty. \end{aligned}$$

Any of these equivalent conditions is equivalent to the function being measurable.

Proof: First note that the first and the third are equivalent. To see this, observe $f^{-1}([d, \infty]) = \bigcap_{n=1}^{\infty} f^{-1}((d - 1/n, \infty])$, and so if the first condition holds, then so does the third. $f^{-1}((d, \infty]) = \bigcup_{n=1}^{\infty} f^{-1}([d + 1/n, \infty])$, and so if the third condition holds, so does the first.

Similarly, the second and fourth conditions are equivalent. Now from the definition of inverse image, $f^{-1}((-\infty, d]) = (f^{-1}((d, \infty]))^C$ so the first and fourth conditions are equivalent. Thus the first four conditions are equivalent and if any of them hold, then for $-\infty < a < b < \infty$, $f^{-1}((a, b)) = f^{-1}((-\infty, b)) \cap f^{-1}((a, \infty]) \in \mathcal{F}$. Finally, if the last condition holds, $f^{-1}([d, \infty]) = (\bigcup_{k=1}^{\infty} f^{-1}((-k + d, \infty)))^C \in \mathcal{F}$ and so the third condition holds. Therefore, all five conditions are equivalent.

Since $(-\infty, \infty]$ is a separable metric space, it follows from Theorem 3.4.2 that every open set U is a countable union of open intervals $U = \bigcup_k I_k$ where I_k is of the form (a, b) or $(a, \infty]$ and, as just shown if any of the equivalent conditions holds, then $f^{-1}(U) = \bigcup_k f^{-1}(I_k) \in \mathcal{F}$. Conversely, if $f^{-1}(U) \in \mathcal{F}$ for any open set $U \in (-\infty, \infty]$, then in particular, $f^{-1}((a, b)) \in \mathcal{F}$ which is one of the equivalent conditions and so all the equivalent conditions hold. ■

Note that if f is continuous and g is measurable, then $f \circ g$ is always measurable. This is because, for U open, $(f \circ g)^{-1}(U) = g^{-1}(f^{-1}(U)) = g^{-1}(\text{open})$ which is measurable.

There is a fundamental theorem about the relationship of simple functions to measurable functions given in the next theorem.

Definition 9.1.5 Let $E \in \mathcal{F}$ for \mathcal{F} a σ algebra. Then

$$\mathcal{X}_E(\omega) \equiv \begin{cases} 1 & \text{if } \omega \in E \\ 0 & \text{if } \omega \notin E \end{cases}$$

This is called the indicator function of the set E . Let $s : (\Omega, \mathcal{F}) \rightarrow \mathbb{R}$. Then s is a simple function if it is of the form

$$s(\omega) = \sum_{i=1}^n c_i \mathcal{X}_{E_i}(\omega)$$

where $E_i \in \mathcal{F}$ and $c_i \in \mathbb{R}$, the E_i being disjoint. Thus simple functions are those which have finitely many values and are measurable. In the next theorem, it will also be assumed that each $c_i \geq 0$.

Each simple function is measurable. This is easily seen as follows. First of all, you can assume the c_i are distinct because if not, you could just replace those E_i which correspond to a single value with their union. Then if you have any open interval (a, b) , $s^{-1}((a, b)) = \bigcup \{E_i : c_i \in (a, b)\}$ and this is measurable because it is the finite union of measurable sets.

Theorem 9.1.6 Let $f \geq 0$ be measurable. Then there exists a sequence of nonnegative simple functions $\{s_n\}$ satisfying

$$0 \leq s_n(\omega) \tag{9.1}$$

$$\cdots s_n(\omega) \leq s_{n+1}(\omega) \cdots$$

$$f(\omega) = \lim_{n \rightarrow \infty} s_n(\omega) \text{ for all } \omega \in \Omega. \tag{9.2}$$

If f is bounded, the convergence is actually uniform. Conversely, if f is nonnegative and is the pointwise limit of such simple functions, then f is measurable.

Proof: Letting $I \equiv \{\omega : f(\omega) = \infty\}$, define

$$t_n(\omega) = \sum_{k=0}^{2^n} \frac{k}{n} \mathcal{X}_{f^{-1}([\frac{k}{n}, \frac{k+1}{n}])}(\omega) + 2^n \mathcal{X}_I(\omega).$$

Then $t_n(\omega) \leq f(\omega)$ for all ω and $\lim_{n \rightarrow \infty} t_n(\omega) = f(\omega)$ for all ω . This is because $t_n(\omega) = 2^n$ for $\omega \in I$ and if $f(\omega) \in [0, \frac{2^n+1}{n})$, then

$$0 \leq f(\omega) - t_n(\omega) \leq \frac{1}{n}. \quad (9.3)$$

Thus whenever $\omega \notin I$, the above inequality will hold for all n large enough. Let

$$s_1 = t_1, s_2 = \max(t_1, t_2), s_3 = \max(t_1, t_2, t_3), \dots$$

Then the sequence $\{s_n\}$ satisfies 9.1-9.2. Also each s_n has finitely many values and is measurable. To see this, note that $s_n^{-1}((a, \infty]) = \cup_{k=1}^n t_k^{-1}((a, \infty]) \in \mathcal{F}$

To verify the last claim, note that in this case the term $2^n \mathcal{X}_I(\omega)$ is not present and for n large enough, $2^n/n$ is larger than all values of f . Therefore, for all n large enough, 9.3 holds for all ω . Thus the convergence is uniform.

The last claim follows right away from Theorem 9.1.2. ■

There is a more general theorem which applies to measurable functions which have values in a separable metric space. In this context, a simple function is one which is of the form $\sum_{k=1}^m x_k \mathcal{X}_{E_k}(\omega)$ where the E_k are disjoint measurable sets and the x_k are in X . I am abusing notation somewhat by using a sum. You can't add in a general metric space. The symbol means the function has value x_k on the set E_k . However, if X were a vector space, this notation would be a nice way to express what is meant.

Theorem 9.1.7 *Let (Ω, \mathcal{F}) be a measurable space and let $f : \Omega \rightarrow X$ where (X, d) is a separable metric space. Then f is a measurable function if and only if there exists a sequence of simple functions, $\{f_n\}$ such that for each $\omega \in \Omega$ and $n \in \mathbb{N}$,*

$$d(f_n(\omega), f(\omega)) \geq d(f_{n+1}(\omega), f(\omega)) \quad (9.4)$$

and

$$\lim_{n \rightarrow \infty} d(f_n(\omega), f(\omega)) = 0. \quad (9.5)$$

Proof: Let $D = \{x_k\}_{k=1}^\infty$ be a countable dense subset of X . First suppose f is measurable. Then since in a metric space every closed set C is the countable intersection of open sets,

$$C = \cap_{k=1}^\infty \{x \in X : \text{dist}(x, C) < 1/k\},$$

it follows $f^{-1}(\text{closed set}) \in \mathcal{F}$. Now let $D_n = \{x_k\}_{k=1}^n$. Let

$$\begin{aligned} A_1 &\equiv \left\{ \omega : d(x_1, f(\omega)) = \min_{k \leq n} d(x_k, f(\omega)) \right\} \\ &= \cap_{k=1}^n \{ \omega : d(x_k, f(\omega)) - d(x_1, f(\omega)) \geq 0 \} \end{aligned}$$

That is, A_1 is those ω such that $f(\omega)$ is approximated best out of D_n by x_1 . Why is this a measurable set? It is because $\omega \rightarrow d(x_k, f(\omega)) - d(x_1, f(\omega))$ is a real valued measurable

function, being the composition of a continuous function, $y \rightarrow d(x_k, y) - d(x_1, y)$ and a measurable function, $\omega \rightarrow f(\omega)$. Next let

$$A_2 \equiv \left\{ \omega \notin A_1 : d(x_2, f(\omega)) = \min_{k \leq n} d(x_k, f(\omega)) \right\}$$

and continue in this manner obtaining disjoint measurable sets, $\{A_k\}_{k=1}^n$ such that for $\omega \in A_k$ the best approximation to $f(\omega)$ from D_n is x_k . Then $f_n(\omega) \equiv \sum_{k=1}^n x_k \chi_{A_k}(\omega)$. Note

$$d(f_{n+1}(\omega), f(\omega)) = \min_{k \leq n+1} d(x_k, f(\omega)) \leq \min_{k \leq n} d(x_k, f(\omega)) = d(f_n(\omega), f(\omega))$$

and so this verifies 9.4. It remains to verify 9.5.

Let $\varepsilon > 0$ be given and pick $\omega \in \Omega$. Then there exists $x_n \in D$ such that $d(x_n, f(\omega)) < \varepsilon$. It follows from the construction that

$$d(f_n(\omega), f(\omega)) \leq d(x_n, f(\omega)) < \varepsilon.$$

This proves the first half.

Conversely, suppose the existence of the sequence of simple functions as described above. Each f_n is a measurable function because $f_n^{-1}(U) = \cup \{A_k : x_k \in U\}$. Therefore, the conclusion that f is measurable follows from Theorem 9.1.2 on Page 237. ■

Another useful observation is that the set where a sequence of measurable functions converges is also a measurable set.

Proposition 9.1.8 *Let $\{f_n\}$ be measurable with values in a complete normed vector space. Let $A \equiv \{\omega : \{f_n(\omega)\} \text{ converges}\}$. Then A is measurable.*

Proof: The set A is the same as the set on which $\{f_n(\omega)\}$ is a Cauchy sequence. This set is

$$\bigcap_{n=1}^{\infty} \bigcup_{m=1}^{\infty} \bigcap_{p,q>m} \left[\|f_p(\omega) - f_q(\omega)\| < \frac{1}{n} \right]$$

which is a measurable set thanks to the measurability of each f_n . ■

9.2 Measures and their Properties

What is meant by a **measure**?

Definition 9.2.1 *Let (Ω, \mathcal{F}) be a measurable space. Here \mathcal{F} is a σ algebra of sets of Ω . Then $\mu : \mathcal{F} \rightarrow [0, \infty]$ is called a measure if whenever $\{F_i\}_{i=1}^{\infty}$ is a sequence of disjoint sets of \mathcal{F} , it follows that*

$$\mu\left(\bigcup_{i=1}^{\infty} F_i\right) = \sum_{i=1}^{\infty} \mu(F_i)$$

Note that the series could equal ∞ . If $\mu(\Omega) < \infty$, then μ is called a finite measure. An important case is when $\mu(\Omega) = 1$ when it is called a probability measure.

Note that $\mu(\emptyset) = \mu(\emptyset \cup \emptyset) = \mu(\emptyset) + \mu(\emptyset)$ and so $\mu(\emptyset) = 0$.

Example 9.2.2 *You could have $\mathcal{P}(\mathbb{N}) = \mathcal{F}$ and you could define $\mu(S)$ to be the number of elements of S . This is called counting measure. It is left as an exercise to show that this is a measure.*

Example 9.2.3 Here is a pathological example. Let Ω be uncountable and \mathcal{F} will be those sets which have the property that either the set is countable or its complement is countable. Let $\mu(E) = 0$ if E is countable and $\mu(E) = 1$ if E is uncountable. It is left as an exercise to show that this is a measure.

Of course the most important measure in this book will be Lebesgue measure which gives the “volume” of a subset of \mathbb{R}^n . However, this requires a lot more work.

Lemma 9.2.4 If μ is a measure and $F_i \in \mathcal{F}$, then $\mu(\cup_{i=1}^{\infty} F_i) \leq \sum_{i=1}^{\infty} \mu(F_i)$. Also if $F_n \in \mathcal{F}$ and $F_n \subseteq F_{n+1}$ for all n , then if $F = \cup_n F_n$,

$$\mu(F) = \lim_{n \rightarrow \infty} \mu(F_n)$$

If $F_n \supseteq F_{n+1}$ for all n , then if $\mu(F_1) < \infty$ and $F = \cap_n F_n$, then

$$\mu(F) = \lim_{n \rightarrow \infty} \mu(F_n)$$

Proof: Let $G_1 = F_1$ and if G_1, \dots, G_n have been chosen disjoint, let $G_{n+1} \equiv F_{n+1} \setminus \cup_{i=1}^n G_i$. Thus the G_i are disjoint. In addition, these are all measurable sets. Now

$$\mu(G_{n+1}) + \mu(F_{n+1} \cap (\cup_{i=1}^n G_i)) = \mu(F_{n+1})$$

and so $\mu(G_n) \leq \mu(F_n)$. Therefore,

$$\mu(\cup_{i=1}^{\infty} G_i) = \sum_i \mu(G_i) \leq \sum_i \mu(F_i).$$

Now consider the increasing sequence of $F_n \in \mathcal{F}$. If $F \subseteq G$ and these are sets of \mathcal{F} , then $\mu(G) = \mu(F) + \mu(G \setminus F)$ so $\mu(G) \geq \mu(F)$. Also $F = \cup_{i=1}^{\infty} (F_{i+1} \setminus F_i) + F_1$. Then $\mu(F) = \sum_{i=1}^{\infty} \mu(F_{i+1} \setminus F_i) + \mu(F_1)$. Now $\mu(F_{i+1} \setminus F_i) + \mu(F_i) = \mu(F_{i+1})$. If any $\mu(F_i) = \infty$, there is nothing to prove. Assume then that these are all finite. Then $\mu(F_{i+1} \setminus F_i) = \mu(F_{i+1}) - \mu(F_i)$ and so

$$\begin{aligned} \mu(F) &= \sum_{i=1}^{\infty} \mu(F_{i+1}) - \mu(F_i) + \mu(F_1) \\ &= \lim_{n \rightarrow \infty} \sum_{i=1}^n \mu(F_{i+1}) - \mu(F_i) + \mu(F_1) = \lim_{n \rightarrow \infty} \mu(F_{n+1}) \end{aligned}$$

Next suppose $\mu(F_1) < \infty$ and $\{F_n\}$ is a decreasing sequence. Then $F_1 \setminus F_n$ is increasing to $F_1 \setminus F$ and so by the first part,

$$\mu(F_1) - \mu(F) = \mu(F_1 \setminus F) = \lim_{n \rightarrow \infty} \mu(F_1 \setminus F_n) = \lim_{n \rightarrow \infty} (\mu(F_1) - \mu(F_n))$$

This is justified because $\mu(F_1 \setminus F_n) + \mu(F_n) = \mu(F_1)$ and all numbers are finite by assumption. Hence $\mu(F) = \lim_{n \rightarrow \infty} \mu(F_n)$. ■

I like to remember this as $E_n \uparrow E \Rightarrow \mu(E_n) \uparrow \mu(E)$ and $E_n \downarrow E \Rightarrow \mu(E_n) \downarrow \mu(E)$ if $\mu(E_1) < \infty$.

There is a monumentally important theorem called the Borel Cantelli lemma. This is next.

Lemma 9.2.5 *If $(\Omega, \mathcal{F}, \mu)$ is a measure space and if $\{E_i\} \subseteq \mathcal{F}$ and $\sum_{i=1}^{\infty} \mu(E_i) < \infty$, then there exists a set N of measure 0 ($\mu(N) = 0$) such that if $\omega \notin N$, then ω is in only finitely many of the E_i .*

Proof: The set of ω in infinitely many E_i is $N \equiv \bigcap_{n=1}^{\infty} \bigcup_{k \geq n} E_k$ because this consists of those ω which are in some E_k for $k \geq n$ for any choice of n . Now $\mu(N) \leq \sum_{k=n}^{\infty} \mu(E_k)$ which is just the tail of a convergent series. Thus, it converges to 0 as $n \rightarrow \infty$. Hence it is less than ε for n large enough. Thus $\mu(N)$ is no more than ε for any $\varepsilon > 0$. ■

9.3 Dynkin's Lemma

Dynkin's lemma is a very useful result. It is used quite a bit in books on probability. It resembles an important result on monotone classes but seems easier to use.

Definition 9.3.1 *Let Ω be a set and let \mathcal{K} be a collection of subsets of Ω . Then \mathcal{K} is called a π system if $\emptyset, \Omega \in \mathcal{K}$ and whenever $A, B \in \mathcal{K}$, it follows $A \cap B \in \mathcal{K}$.*

The following is the fundamental lemma which shows these π systems are useful. This is due to Dynkin.

Lemma 9.3.2 *Let \mathcal{K} be a π system of subsets of Ω , a set. Also let \mathcal{G} be a collection of subsets of Ω which satisfies the following three properties.*

1. $\mathcal{K} \subseteq \mathcal{G}$
2. If $A \in \mathcal{G}$, then $A^C \in \mathcal{G}$
3. If $\{A_i\}_{i=1}^{\infty}$ is a sequence of disjoint sets from \mathcal{G} then $\bigcup_{i=1}^{\infty} A_i \in \mathcal{G}$.

Then $\mathcal{G} \supseteq \sigma(\mathcal{K})$, where $\sigma(\mathcal{K})$ is the smallest σ algebra which contains \mathcal{K} .

Proof: First note that if

$$\mathcal{H} \equiv \{\mathcal{G} : \text{1 - 3 all hold}\}$$

then $\bigcap \mathcal{H}$ yields a collection of sets which also satisfies 1 - 3. Therefore, I will assume in the argument that \mathcal{G} is the smallest collection satisfying 1 - 3. Let $A \in \mathcal{K}$ and define

$$\mathcal{G}_A \equiv \{B \in \mathcal{G} : A \cap B \in \mathcal{G}\}.$$

I want to show \mathcal{G}_A satisfies 1 - 3 because then it must equal \mathcal{G} since \mathcal{G} is the smallest collection of subsets of Ω which satisfies 1 - 3. This will give the conclusion that for $A \in \mathcal{K}$ and $B \in \mathcal{G}$, $A \cap B \in \mathcal{G}$. This information will then be used to show that if $A, B \in \mathcal{G}$ then $A \cap B \in \mathcal{G}$. From this it will follow very easily that \mathcal{G} is a σ algebra which will imply it contains $\sigma(\mathcal{K})$. Now here are the details of the argument.

Since \mathcal{K} is given to be a π system contained in \mathcal{G} , $\mathcal{K} \subseteq \mathcal{G}_A$. Indeed, if $C \in \mathcal{K}$ then $A \cap C \in \mathcal{K} \subseteq \mathcal{G}$ so $C \in \mathcal{G}_A$. Property 3 is obvious because if $\{B_i\}$ is a sequence of disjoint sets in \mathcal{G}_A , then

$$A \cap \bigcup_{i=1}^{\infty} B_i = \bigcup_{i=1}^{\infty} A \cap B_i \in \mathcal{G}$$

because $A \cap B_i \in \mathcal{G}$ and the property 3 of \mathcal{G} .

It remains to verify Property 2 so let $B \in \mathcal{G}_A$. I need to verify that $B^C \in \mathcal{G}_A$. In other words, I need to show that $A \cap B^C \in \mathcal{G}$. However,

$$A \cap B^C = (A^C \cup (A \cap B))^C \in \mathcal{G}$$

Here is why. Since $B \in \mathcal{G}_A$, $A \cap B \in \mathcal{G}$ and since $A \in \mathcal{K} \subseteq \mathcal{G}$ it follows $A^C \in \mathcal{G}$ by assumption 2. It follows from assumption 3 the union of the disjoint sets, A^C and $(A \cap B)$ is in \mathcal{G} and then from 2 the complement of their union is in \mathcal{G} . Thus \mathcal{G}_A satisfies 1 - 3 and this implies, since \mathcal{G} is the smallest such, that $\mathcal{G}_A \supseteq \mathcal{G}$. However, \mathcal{G}_A is constructed as a subset of \mathcal{G} . This proves that for every $B \in \mathcal{G}$ and $A \in \mathcal{K}$, $A \cap B \in \mathcal{G}$. Now pick $B \in \mathcal{G}$ and consider

$$\mathcal{G}_B \equiv \{A \in \mathcal{G} : A \cap B \in \mathcal{G}\}.$$

I just proved $\mathcal{K} \subseteq \mathcal{G}_B$. The other arguments are identical to show \mathcal{G}_B satisfies 1 - 3 and is therefore equal to \mathcal{G} . This shows that whenever $A, B \in \mathcal{G}$ it follows $A \cap B \in \mathcal{G}$.

This implies \mathcal{G} is a σ algebra. To show this, all that is left is to verify \mathcal{G} is closed under countable unions because then it follows \mathcal{G} is a σ algebra. Let $\{A_i\} \subseteq \mathcal{G}$. Then let $A'_1 = A_1$ and

$$A'_{n+1} \equiv A_{n+1} \setminus (\cup_{i=1}^n A_i) = A_{n+1} \cap (\cap_{i=1}^n A_i^C) = \cap_{i=1}^n (A_{n+1} \cap A_i^C) \in \mathcal{G}$$

because finite intersections of sets of \mathcal{G} are in \mathcal{G} . Since the A'_i are disjoint, it follows $\cup_{i=1}^\infty A_i = \cup_{i=1}^\infty A'_i \in \mathcal{G}$. Therefore, $\mathcal{G} \supseteq \sigma(\mathcal{K})$. ■

Corollary 9.3.3 *Given 2, closed with respect to complements, the condition that \mathcal{G} is closed with respect to countable disjoint unions is equivalent to \mathcal{G} the condition that \mathcal{G} is closed with respect to countable intersections.*

Proof: \Rightarrow Consider $\cap_k E_k$ where $E_k \in \mathcal{G}$. Then $\cap_k E_k = (\cup_k E_k^C)^C$. Now the E_k^C are not necessarily disjoint, but each is in \mathcal{G} and so one can use the scheme of the last part of the proof of Lemma 9.3.2 to reduce to this case and conclude $\cup_k E_k^C \in \mathcal{G}$. Then the countable intersection is just the complement of this last set.

\Leftarrow Suppose the countable intersection of sets of \mathcal{G} is in \mathcal{G} and consider a countable union $\cup_k E_k$ of sets of \mathcal{G} . Then $\cup_k E_k = (\cap_k E_k^C)^C \in \mathcal{G}$. ■

9.4 Outer Measures

There is also something called an outer measure which is defined on the set of all subsets.

Definition 9.4.1 *Let Ω be a nonempty set and let $\lambda : \mathcal{P}(\Omega) \rightarrow [0, \infty)$ satisfy the following:*

1. $\lambda(\emptyset) = 0$
2. If $A \subseteq B$, then $\lambda(A) \leq \lambda(B)$
3. $\lambda(\cup_{i=1}^\infty E_i) \leq \sum_{i=1}^\infty \lambda(E_i)$

Then λ is called an outer measure.

Every measure determines an outer measure. For example, suppose that μ is a measure on \mathcal{F} a σ algebra of subsets of Ω . Then define

$$\bar{\mu}(S) \equiv \inf \{ \mu(E) : E \supseteq S, E \in \mathcal{F} \}.$$

This is easily seen to be an outer measure. Also, we have the following Proposition.

Proposition 9.4.2 *Let μ be a measure defined on a σ algebra of subsets \mathcal{F} of Ω as just described. Then $\bar{\mu}$ as defined above, is an outer measure and also, if $E \in \mathcal{F}$, then $\bar{\mu}(E) = \mu(E)$.*

Proof: The first two properties of an outer measure are obvious. What of the third? If any $\bar{\mu}(E_i) = \infty$, then there is nothing to show so suppose each of these is finite. Let $F_i \supseteq E_i$ such that $F_i \in \mathcal{F}$ and $\bar{\mu}(E_i) + \frac{\varepsilon}{2^i} > \mu(F_i)$. Then

$$\bar{\mu}(\cup_{i=1}^{\infty} E_i) \leq \mu(\cup_{i=1}^{\infty} F_i) \leq \sum_{i=1}^{\infty} \mu(F_i) < \sum_{i=1}^{\infty} \left(\bar{\mu}(E_i) + \frac{\varepsilon}{2^i} \right) = \sum_{i=1}^{\infty} \bar{\mu}(E_i) + \varepsilon$$

Since ε is arbitrary, this establishes the third condition. Finally, if $E \in \mathcal{F}$, then by definition, $\bar{\mu}(E) \leq \mu(E)$ because $E \supseteq E$. Also, $\mu(E) \leq \mu(F)$ for all $F \in \mathcal{F}$ such that $F \supseteq E$. It follows that $\mu(E)$ is a lower bound of all such $\mu(F)$ and so $\bar{\mu}(E) \geq \mu(E)$. ■

9.5 Measures From Outer Measures

Theorem 9.7.4 describes an outer measure on $\mathcal{P}(\mathbb{R})$. There is a general procedure for constructing a σ algebra and a measure from an outer measure which is due to Caratheodory about 1918.

Thus, when you have a measure on (Ω, \mathcal{F}) , you can obtain an outer measure on $(\Omega, \mathcal{P}(\Omega))$ from this measure as in Proposition 9.4.2, and if you have an outer measure on $(\Omega, \mathcal{P}(\Omega))$, this will define a σ algebra \mathcal{F} and a measure on (Ω, \mathcal{F}) . This last assertion is the topic of this section.

Definition 9.5.1 *Let Ω be a nonempty set and let $\mu : \mathcal{P}(\Omega) \rightarrow [0, \infty]$ be an outer measure. For $E \subseteq \Omega$, E is μ measurable if for all $S \subseteq \Omega$,*

$$\mu(S) = \mu(S \setminus E) + \mu(S \cap E). \quad (9.6)$$

To help in remembering 9.6, think of a measurable set E , as a process which divides a given set into two pieces, the part in E and the part not in E as in 9.6. In the Bible, there are several incidents recorded in which a process of division resulted in more stuff than was originally present.¹ Measurable sets are exactly those which are incapable of such a miracle. With an outer measure, it is always the case that $\mu(S) \leq \mu(S \setminus E) + \mu(S \cap E)$. The set is measurable, when equality is always obtained for any choice of $S \in \mathcal{P}(\Omega)$. You might think of the measurable sets as the non-miraculous sets. The idea is to show that these sets form a σ algebra on which the outer measure μ is a measure.

First here is a definition and a lemma.

¹ 1 Kings 17, 2 Kings 4, Mathew 14, and Mathew 15 all contain such descriptions. The stuff involved was either oil, bread, flour or fish. In mathematics such things have also been done with sets. In the book by Bruckner Bruckner and Thompson there is an interesting discussion of the Banach Tarski paradox which says it is possible to divide a ball in \mathbb{R}^3 into five disjoint pieces and assemble the pieces to form two disjoint balls of the same size as the first. The details can be found in: The Banach Tarski Paradox by Wagon, Cambridge University press. 1985. It is known that all such examples must involve the axiom of choice.

Definition 9.5.2 $(\mu \lfloor S)(A) \equiv \mu(S \cap A)$ for all $A \subseteq \Omega$. Thus $\mu \lfloor S$ is the name of a new outer measure, called μ restricted to S .

The next lemma indicates that the property of measurability is not lost by considering this restricted measure.

Lemma 9.5.3 If A is μ measurable, then for any S , A is $\mu \lfloor S$ measurable.

Proof: Suppose A is μ measurable. It is desired to show that for all $T \subseteq \Omega$,

$$(\mu \lfloor S)(T) = (\mu \lfloor S)(T \cap A) + (\mu \lfloor S)(T \setminus A).$$

Thus it is desired to show

$$\mu(S \cap T) = \mu(T \cap A \cap S) + \mu(T \cap S \cap A^C). \quad (9.7)$$

But 9.7 holds because A is μ measurable. Apply Definition 9.5.1 to $S \cap T$ instead of S . ■

If A is $\mu \lfloor S$ measurable, it does not follow that A is μ measurable. Indeed, if you believe in the existence of non measurable sets which is discussed later, you could let $A = S$ for such a μ non measurable set and verify that S is $\mu \lfloor S$ measurable.

The next theorem is the main result on outer measures which shows that starting with an outer measure you can obtain a measure.

Theorem 9.5.4 Let Ω be a set and let μ be an outer measure on $\mathcal{P}(\Omega)$. The collection of μ measurable sets \mathcal{S} , forms a σ algebra and

$$\text{If } F_i \in \mathcal{S}, F_i \cap F_j = \emptyset, \text{ then } \mu(\cup_{i=1}^{\infty} F_i) = \sum_{i=1}^{\infty} \mu(F_i). \quad (9.8)$$

If $\cdots F_n \subseteq F_{n+1} \subseteq \cdots$, then if $F = \cup_{n=1}^{\infty} F_n$ and $F_n \in \mathcal{S}$, it follows that

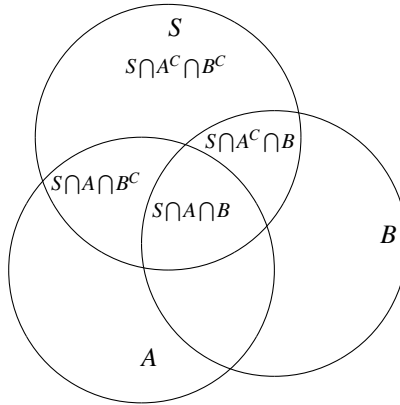
$$\mu(F) = \lim_{n \rightarrow \infty} \mu(F_n). \quad (9.9)$$

If $\cdots F_n \supseteq F_{n+1} \supseteq \cdots$, and if $F = \cap_{n=1}^{\infty} F_n$ for $F_n \in \mathcal{S}$ then if $\mu(F_1) < \infty$,

$$\mu(F) = \lim_{n \rightarrow \infty} \mu(F_n). \quad (9.10)$$

This measure space is also complete which means that if $\mu(F) = 0$ for some $F \in \mathcal{S}$ then if $G \subseteq F$, it follows $G \in \mathcal{S}$ also.

Proof: First note that \emptyset and Ω are obviously in \mathcal{S} . Now suppose $A, B \in \mathcal{S}$. I will show $A \setminus B \equiv A \cap B^C$ is in \mathcal{S} . To do so, consider the following picture.



It is required to show that $\mu(S) = \mu(S \setminus (A \setminus B)) + \mu(S \cap (A \setminus B))$. First consider $S \setminus (A \setminus B)$. From the picture, it equals

$$(S \cap A^C \cap B^C) \cup (S \cap A \cap B) \cup (S \cap A^C \cap B)$$

Therefore, $\mu(S) \leq \mu(S \setminus (A \setminus B)) + \mu(S \cap (A \setminus B))$

$$\begin{aligned} &\leq \mu(S \cap A^C \cap B^C) + \mu(S \cap A \cap B) + \mu(S \cap A^C \cap B) + \mu(S \cap (A \setminus B)) \\ &= \mu(S \cap A^C \cap B^C) + \mu(S \cap A \cap B) + \mu(S \cap A^C \cap B) + \mu(S \cap A \cap B^C) \\ &= \mu(S \cap A^C \cap B^C) + \mu(S \cap A \cap B^C) + \mu(S \cap A \cap B) + \mu(S \cap A^C \cap B) \\ &= \mu(S \cap B^C) + \mu(S \cap B) = \mu(S) \end{aligned}$$

and so this shows that $A \setminus B \in \mathcal{S}$ whenever $A, B \in \mathcal{S}$.

Since $\Omega \in \mathcal{S}$, this shows that $A \in \mathcal{S}$ if and only if $A^C \in \mathcal{S}$. Now if $A, B \in \mathcal{S}$, $A \cup B = (A^C \cap B^C)^C = (A^C \setminus B)^C \in \mathcal{S}$. By induction, if $A_1, \dots, A_n \in \mathcal{S}$, then so is $\cup_{i=1}^n A_i$. If $A, B \in \mathcal{S}$, with $A \cap B = \emptyset$,

$$\mu(A \cup B) = \mu((A \cup B) \cap A) + \mu((A \cup B) \setminus A) = \mu(A) + \mu(B).$$

By induction, if $A_i \cap A_j = \emptyset$ and $A_i \in \mathcal{S}$,

$$\mu(\cup_{i=1}^n A_i) = \sum_{i=1}^n \mu(A_i). \quad (9.11)$$

Now let $A = \cup_{i=1}^\infty A_i$ where $A_i \cap A_j = \emptyset$ for $i \neq j$. $\sum_{i=1}^\infty \mu(A_i) \geq \mu(A) \geq \mu(\cup_{i=1}^n A_i) = \sum_{i=1}^n \mu(A_i)$. Since this holds for all n , you can take the limit as $n \rightarrow \infty$ and conclude, $\sum_{i=1}^\infty \mu(A_i) = \mu(A)$ which establishes 9.8.

Consider part 9.9. Without loss of generality $\mu(F_k) < \infty$ for all k since otherwise there is nothing to show. Suppose $\{F_k\}$ is an increasing sequence of sets of \mathcal{S} . Then letting $F_0 \equiv \emptyset$, $\{F_{k+1} \setminus F_k\}_{k=0}^\infty$ is a sequence of disjoint sets of \mathcal{S} since it was shown above that the difference of two sets of \mathcal{S} is in \mathcal{S} . Also note that from 9.11

$$\mu(F_{k+1} \setminus F_k) + \mu(F_k) = \mu(F_{k+1})$$

and so if $\mu(F_k) < \infty$, then

$$\mu(F_{k+1} \setminus F_k) = \mu(F_{k+1}) - \mu(F_k).$$

Therefore, letting $F \equiv \cup_{k=1}^\infty F_k$ which also equals $\cup_{k=1}^\infty (F_{k+1} \setminus F_k)$, it follows from part 9.8 just shown that

$$\begin{aligned} \mu(F) &= \sum_{k=0}^\infty \mu(F_{k+1} \setminus F_k) = \lim_{n \rightarrow \infty} \sum_{k=0}^n \mu(F_{k+1} \setminus F_k) \\ &= \lim_{n \rightarrow \infty} \sum_{k=0}^n \mu(F_{k+1}) - \mu(F_k) = \lim_{n \rightarrow \infty} \mu(F_{n+1}). \end{aligned}$$

In order to establish 9.10, let the F_n be as given there. Then, since $(F_1 \setminus F_n)$ increases to $(F_1 \setminus F)$, 9.9 implies

$$\lim_{n \rightarrow \infty} (\mu(F_1) - \mu(F_n)) = \lim_{n \rightarrow \infty} \mu(F_1 \setminus F_n) = \mu(F_1 \setminus F).$$

The problem is, I don't know $F \in \mathcal{S}$ and so it is not clear that $\mu(F_1 \setminus F) = \mu(F_1) - \mu(F)$. However, $\mu(F_1 \setminus F) + \mu(F) \geq \mu(F_1)$ and so $\mu(F_1 \setminus F) \geq \mu(F_1) - \mu(F)$. Hence

$$\lim_{n \rightarrow \infty} (\mu(F_1) - \mu(F_n)) = \mu(F_1 \setminus F) \geq \mu(F_1) - \mu(F)$$

which implies $\lim_{n \rightarrow \infty} \mu(F_n) \leq \mu(F)$. But since $F \subseteq F_n$, $\mu(F) \leq \lim_{n \rightarrow \infty} \mu(F_n)$ and this establishes 9.10. Note that it was assumed $\mu(F_1) < \infty$ because $\mu(F_1)$ was subtracted from both sides.

It remains to show \mathcal{S} is closed under countable unions. Recall that if $A \in \mathcal{S}$, then $A^C \in \mathcal{S}$ and \mathcal{S} is closed under finite unions. Let $A_i \in \mathcal{S}$, $A = \cup_{i=1}^{\infty} A_i$, $B_n = \cup_{i=1}^n A_i$. Then

$$\begin{aligned} \mu(S) &= \mu(S \cap B_n) + \mu(S \setminus B_n) \\ &= (\mu \lfloor S)(B_n) + (\mu \lfloor S)(B_n^C). \end{aligned} \quad (9.12)$$

By Lemma 9.5.3 B_n is $(\mu \lfloor S)$ measurable and so is B_n^C . I want to show $\mu(S) \geq \mu(S \setminus A) + \mu(S \cap A)$. If $\mu(S) = \infty$, there is nothing to prove. Assume $\mu(S) < \infty$. Then apply Parts 9.10 and 9.9 to the outer measure $\mu \lfloor S$ in 9.12 and let $n \rightarrow \infty$. Thus $B_n \uparrow A$, $B_n^C \downarrow A^C$ and this yields $\mu(S) = (\mu \lfloor S)(A) + (\mu \lfloor S)(A^C) = \mu(S \cap A) + \mu(S \setminus A)$.

Therefore $A \in \mathcal{S}$ and this proves Parts 9.8, 9.9, and 9.10.

It only remains to verify the assertion about completeness. Letting G and F be as described above, let $S \subseteq \Omega$. I need to verify $\mu(S) \geq \mu(S \cap G) + \mu(S \setminus G)$. However,

$$\begin{aligned} \mu(S \cap G) + \mu(S \setminus G) &\leq \mu(S \cap F) + \mu(S \setminus F) + \mu(F \setminus G) \\ &= \mu(S \cap F) + \mu(S \setminus F) = \mu(S) \end{aligned}$$

because by assumption, $\mu(F \setminus G) \leq \mu(F) = 0$. ■

Corollary 9.5.5 *Completeness is the same as saying that if $(E \setminus E') \cup (E' \setminus E) \subseteq N \in \mathcal{F}$ and $\mu(N) = 0$, then if $E \in \mathcal{F}$, it follows that $E' \in \mathcal{F}$ also.*

Proof: If the new condition holds, then suppose $G \subseteq F$ where $\mu(F) = 0$, $F \in \mathcal{F}$. Then $\overbrace{(G \setminus F)}^{=\emptyset} \cup (F \setminus G) \subseteq F$ and $\mu(F)$ is given to equal 0. Therefore, $G \in \mathcal{F}$.

Now suppose the earlier version of completeness and let

$$(E \setminus E') \cup (E' \setminus E) \subseteq N \in \mathcal{F}$$

where $\mu(N) = 0$ and $E \in \mathcal{F}$. Then we know $(E \setminus E'), (E' \setminus E) \in \mathcal{F}$ and all have measure zero. It follows $E \setminus (E \setminus E') = E \cap E' \in \mathcal{F}$. Hence

$$E' = (E \cap E') \cup (E' \setminus E) \in \mathcal{F} \quad \blacksquare$$

9.6 Measurable Sets Include Borel Sets?

If you have an outer measure, it determines a measure. This section gives a very convenient criterion which allows you to conclude right away that the measure is a Borel measure.

Theorem 9.6.1 *Let μ be an outer measure on the subsets of (X, d) , a metric space. If $\mu(A \cup B) = \mu(A) + \mu(B)$ whenever $\text{dist}(A, B) > 0$, then the σ algebra of measurable sets \mathcal{S} contains the Borel sets.*

Proof: It suffices to show that closed sets are in \mathcal{S} , the σ -algebra of measurable sets, because then the open sets are also in \mathcal{S} and consequently \mathcal{S} contains the Borel sets. Let K be closed and let S be a subset of Ω . Is $\mu(S) \geq \mu(S \cap K) + \mu(S \setminus K)$? It suffices to assume $\mu(S) < \infty$. Let $K_n \equiv \{x : \text{dist}(x, K) \leq \frac{1}{n}\}$. By Lemma 3.12.1 on Page 91, $x \rightarrow \text{dist}(x, K)$ is continuous and so K_n is a closed set having K as a subset. That in K_n^C is at a positive distance from K . By the assumption of the theorem,

$$\mu(S) \geq \mu((S \cap K) \cup (S \setminus K_n)) = \mu(S \cap K) + \mu(S \setminus K_n) \quad (9.13)$$

Now

$$\mu(S \setminus K_n) \leq \mu(S \setminus K) \leq \mu(S \setminus K_n) + \mu((K_n \setminus K) \cap S). \quad (9.14)$$

If $\lim_{n \rightarrow \infty} \mu((K_n \setminus K) \cap S) = 0$ then the theorem will be proved because this limit along with 9.14 implies $\lim_{n \rightarrow \infty} \mu(S \setminus K_n) = \mu(S \setminus K)$ and then taking a limit in 9.13, $\mu(S) \geq \mu(S \cap K) + \mu(S \setminus K)$ as desired. Therefore, it suffices to establish this limit.

Since K is closed, a point, $x \notin K$ must be at a positive distance from K and so

$$K_n \setminus K = \bigcup_{k=n}^{\infty} K_k \setminus K_{k+1}.$$

Therefore

$$\mu(S \cap (K_n \setminus K)) \leq \sum_{k=n}^{\infty} \mu(S \cap (K_k \setminus K_{k+1})). \quad (9.15)$$

If

$$\sum_{k=1}^{\infty} \mu(S \cap (K_k \setminus K_{k+1})) < \infty, \quad (9.16)$$

then $\mu(S \cap (K_n \setminus K)) \rightarrow 0$ because it is dominated by the tail of a convergent series so it suffices to show 9.16.

$$\begin{aligned} \sum_{k=1}^M \mu(S \cap (K_k \setminus K_{k+1})) &= \\ \sum_{k \text{ even}, k \leq M} \mu(S \cap (K_k \setminus K_{k+1})) &+ \sum_{k \text{ odd}, k \leq M} \mu(S \cap (K_k \setminus K_{k+1})). \end{aligned} \quad (9.17)$$

By the construction, the distance between any pair of sets, $S \cap (K_k \setminus K_{k+1})$ for different even values of k is positive and the distance between any pair of sets, $S \cap (K_k \setminus K_{k+1})$ for different odd values of k is positive. Therefore,

$$\begin{aligned} \sum_{k \text{ even}, k \leq M} \mu(S \cap (K_k \setminus K_{k+1})) &+ \sum_{k \text{ odd}, k \leq M} \mu(S \cap (K_k \setminus K_{k+1})) \leq \\ \mu \left(\bigcup_{k \text{ even}, k \leq M} (S \cap (K_k \setminus K_{k+1})) \right) &+ \mu \left(\bigcup_{k \text{ odd}, k \leq M} (S \cap (K_k \setminus K_{k+1})) \right) \\ &\leq \mu(S) + \mu(S) = 2\mu(S) \end{aligned}$$

and so for all M , $\sum_{k=1}^M \mu(S \cap (K_k \setminus K_{k+1})) \leq 2\mu(S)$ showing 9.16. ■

9.7 An Outer Measure on $\mathcal{P}(\mathbb{R})$

A measure on \mathbb{R} is like length. I will present something more general than length because it is no trouble to do so and the generalization is useful in many areas of mathematics such as probability.

Definition 9.7.1 *The following definition is important.*

$$F(x+) \equiv \lim_{y \rightarrow x+} F(y), \quad F(x-) \equiv \lim_{y \rightarrow x-} F(y)$$

Thus one of these is the limit from the left and the other is the limit from the right.

In probability, one often has $F(x) \geq 0$, F is increasing, and $F(x+) = F(x)$. This is the case where F is a probability distribution function. In this case, $F(x) \equiv P(X \leq x)$ where X is a random variable. In this case, $\lim_{x \rightarrow \infty} F(x) = 1$ but we are considering more general functions than this including the simple example where $F(x) = x$. This last example will end up giving Lebesgue measure on \mathbb{R} . Recall the following definition.

Definition 9.7.2 $\mathcal{P}(S)$ denotes the set of all subsets of S .

Also recall

Definition 9.7.3 *For two sets, A, B in a metric space,*

$$\text{dist}(A, B) \equiv \inf \{d(x, y) : x \in A, y \in B\}.$$

Theorem 9.7.4 *Let F be an increasing function defined on \mathbb{R} . This will be called an integrator function. There exists a function $\mu : \mathcal{P}(\mathbb{R}) \rightarrow [0, \infty]$ which satisfies the following properties.*

1. *If $A \subseteq B$, then $0 \leq \mu(A) \leq \mu(B)$, $\mu(\emptyset) = 0$.*
2. *$\mu(\cup_{k=1}^{\infty} A_k) \leq \sum_{k=1}^{\infty} \mu(A_k)$*
3. *$\mu([a, b]) = F(b+) - F(a-)$,*
4. *$\mu((a, b)) = F(b-) - F(a+)$*
5. *$\mu((a, b]) = F(b+) - F(a+)$*
6. *$\mu([a, b)) = F(b-) - F(a-)$.*
7. *If $\text{dist}(A, B) = \delta > 0$, then $\mu(A \cup B) = \mu(A) + \mu(B)$.*

Then the σ algebra of μ measurable sets \mathcal{F} contains the Borel sets. This measure is called Lebesgue Stieltjes measure.

Proof: First it is necessary to define the function μ . This is contained in the following definition.

Definition 9.7.5 *For $A \subseteq \mathbb{R}$,*

$$\mu(A) = \inf \left\{ \sum_{i=1}^{\infty} (F(b_i-) - F(a_i+)) : A \subseteq \cup_{i=1}^{\infty} (a_i, b_i) \right\}$$

In words, you look at all coverings of A with open intervals. For each of these open coverings, you add the “lengths” of the individual open intervals and you take the infimum of all such numbers obtained.

Then 1.) is obvious because if a countable collection of open intervals covers B , then it also covers A . Thus the set of numbers obtained for B is smaller than the set of numbers for A . Why is $\mu(\emptyset) = 0$? Pick a point of continuity of F . Such points exist because F is increasing and so it has only countably many points of discontinuity. Let a be this point. Then $\emptyset \subseteq (a - \delta, a + \delta)$ and so $\mu(\emptyset) \leq F(a + \delta) - F(a - \delta)$ for every $\delta > 0$. Letting $\delta \rightarrow 0$, it follows that $\mu(\emptyset) = 0$.

Consider 2.). If any $\mu(A_i) = \infty$, there is nothing to prove. The assertion simply is $\infty \leq \infty$. Assume then that $\mu(A_i) < \infty$ for all i . Then for each $m \in \mathbb{N}$ there exists a countable set of open intervals, $\{(a_i^m, b_i^m)\}_{i=1}^\infty$ such that

$$\mu(A_m) + \frac{\varepsilon}{2^m} > \sum_{i=1}^\infty (F(b_i^m -) - F(a_i^m +)).$$

Then using Theorem 2.5.4 on Page 65,

$$\begin{aligned} \mu(\cup_{m=1}^\infty A_m) &\leq \sum_{i,m} (F(b_i^m -) - F(a_i^m +)) \\ &= \sum_{m=1}^\infty \sum_{i=1}^\infty (F(b_i^m -) - F(a_i^m +)) \leq \sum_{m=1}^\infty \mu(A_m) + \frac{\varepsilon}{2^m} = \sum_{m=1}^\infty \mu(A_m) + \varepsilon, \end{aligned}$$

and since ε is arbitrary, this establishes 2.).

Next consider 3.). By definition, there exists a sequence of open intervals, $\{(a_i, b_i)\}_{i=1}^\infty$ whose union contains $[a, b]$ such that

$$\mu([a, b]) + \varepsilon \geq \sum_{i=1}^\infty (F(b_i -) - F(a_i +)).$$

By Theorem 4.4.8, finitely many of these open intervals also cover $[a, b]$. It follows there exist finitely many of these intervals, denoted as $\{(a_i, b_i)\}_{i=1}^n$, which overlap, such that $a \in (a_1, b_1)$, $b_1 \in (a_2, b_2)$, \dots , $b \in (a_n, b_n)$. Therefore, $\mu([a, b]) \leq \sum_{i=1}^n (F(b_i -) - F(a_i +))$. It follows

$$\begin{aligned} \sum_{i=1}^n (F(b_i -) - F(a_i +)) &\geq \mu([a, b]) \geq \sum_{i=1}^n (F(b_i -) - F(a_i +)) - \varepsilon \\ &\geq F(b+) - F(a-) - \varepsilon \end{aligned}$$

Therefore, $F(b + \delta) - F(a - \delta) \geq \mu([a, b]) \geq F(b+) - F(a-) - \varepsilon$. Letting $\delta \rightarrow 0$,

$$F(b+) - F(a-) \geq \mu([a, b]) \geq F(b+) - F(a-) - \varepsilon$$

Since ε is arbitrary, this shows $\mu([a, b]) = F(b+) - F(a-)$. This establishes 3.).

Consider 4.). For small $\delta > 0$, $\mu([a + \delta, b - \delta]) \leq \mu((a, b))$. Therefore, from 3.) and the definition of μ ,

$$\begin{aligned} F((b - \delta)) - F((a + \delta)) &\leq F((b - \delta) +) - F((a + \delta) -) \\ &= \mu([a + \delta, b - \delta]) \leq \mu((a, b)) \leq F(b-) - F(a+) \end{aligned}$$

the last inequality from the definition. Now letting δ decrease to 0 it follows $F(b-) - F(a+) \leq \mu((a, b)) \leq F(b-) - F(a+)$. This shows 4.)

Consider 5.). From 3.) and 4.), for small $\delta > 0$,

$$\begin{aligned} F(b+) - F((a + \delta)) &\leq F(b+) - F((a + \delta) -) \\ &= \mu([a + \delta, b]) \leq \mu((a, b]) \leq \mu((a, b + \delta)) \\ &= F((b + \delta) -) - F(a+) \leq F(b + \delta) - F(a+). \end{aligned}$$

Now let δ converge to 0 from above to obtain $F(b+) - F(a+) = \mu((a, b])$. This establishes 5.) and 6.) is entirely similar to 5.).

Finally, consider 7.). Let

$$V \equiv \bigcup \left\{ B\left(x, \frac{\delta}{10}\right) : x \in A \cup B \right\}.$$

Let $A \cup B \subseteq \bigcup_{i=1}^{\infty} (a_i, b_i)$ where

$$\mu(A \cup B) + \varepsilon > \sum_i F(b_i-) - F(a_i+)$$

Then, taking the intersection of each of these intervals with V , it can be assumed that all of the intervals are contained in V since such an intersection will only strengthen the above inequality. Now refer to V as the union of these intervals, none of which can intersect both A and B . Thus V consists of disjoint open sets, one containing A consisting of the intervals which intersect A , U_A and the other consisting of those which intersect B , U_B . Let \mathcal{I}_A denote the intervals which intersect A and let \mathcal{I}_B denote the remaining intervals. Also let $\Delta((a_i, b_i)) \equiv F(b_i-) - F(a_i+)$. Then from the above,

$$\mu(A \cup B) + \varepsilon > \sum_{I \in \mathcal{I}_A} \Delta(I) + \sum_{I \in \mathcal{I}_B} \Delta(I) \geq \mu(A) + \mu(B) \geq \mu(A \cup B)$$

Since $\varepsilon > 0$ is arbitrary, this shows 7.). That \mathcal{F} contains the Borel sets follows from 7.) also. ■

We have just shown that μ is an outer measure on $\mathcal{P}(\mathbb{R})$. Unlike what was presented earlier, this outer measure did not begin with a measure.

9.8 Measures and Regularity

It is often the case that Ω is not just a set. In particular, it is often the case that Ω is some sort of topological space, often a metric space. In this case, it is usually if not always the case that the open sets will be in the σ algebra of measurable sets. This leads to the following definition.

Definition 9.8.1 *A Polish space is a complete separable metric space. For a Polish space E or more generally a metric space or even a general topological space, $\mathcal{B}(E)$ denotes the **Borel sets** of E . This is defined to be the smallest σ algebra which contains the open sets. Thus it contains all open sets and closed sets and compact sets and many others.*

For example, \mathbb{R} is a Polish space as is any separable Banach space. **Amazing things** can be said about finite measures on the Borel sets of a Polish space. First the case of a finite measure on a metric space will be considered.

It is best to not attempt to describe a generic Borel set. Always work with the definition that it is the smallest σ algebra containing the open sets. Attempts to give an explicit description of a “typical” Borel set tend to lead nowhere because there are so many things which can be done. You can take countable unions and complements and then countable intersections of what you get and then another countable union followed by complements and on and on. You just can’t get a good useable description in this way. However, it is easy to see that something like $\left(\bigcap_{i=1}^{\infty} \bigcup_{j=i}^{\infty} E_j\right)^C$ is a Borel set if the E_j are. This is useful. This said, you can look at Hewitt and Stromberg [26] in their discussion of why there are more Lebesgue measurable sets than Borel measurable sets to see the kind of technicalities which result by describing Borel sets. This is an extremely significant result based on describing Borel sets, so it can be done.

Definition 9.8.2 A measure μ defined on a σ algebra \mathcal{F} which includes $\mathcal{B}(E)$ will be called inner regular on \mathcal{F} if for all $F \in \mathcal{F}$,

$$\mu(F) = \sup\{\mu(K) : K \subseteq F \text{ and } K \text{ is closed}\} \quad (9.18)$$

A measure, μ defined on \mathcal{F} will be called outer regular on \mathcal{F} if for all $F \in \mathcal{F}$,

$$\mu(F) = \inf\{\mu(V) : V \supseteq F \text{ and } V \text{ is open}\} \quad (9.19)$$

When a measure is both inner and outer regular, it is called regular. Actually, it is more useful and likely more standard to refer to μ being inner regular as

$$\mu(F) = \sup\{\mu(K) : K \subseteq F \text{ and } K \text{ is compact}\} \quad (9.20)$$

Thus the word “closed” is replaced with “compact”. A complete measure defined on a σ algebra \mathcal{F} which includes the Borel sets which is finite on compact sets and also satisfies 9.19 and 9.20 for each $F \in \mathcal{F}$ is called a Radon measure. A G_δ set, pronounced as G delta is the countable intersection of open sets. An F_σ set, pronounced F sigma is the countable union of closed sets.

In every case which has ever been of interest to me, the measure has been σ finite.

Definition 9.8.3 If (X, \mathcal{F}, μ) is a measure space, it is called σ finite if there are $X_n \in \mathcal{F}$ with $\bigcup_n X_n = X$ and $\mu(X_n) < \infty$.

For finite measures, defined on the Borel sets of a metric space X , $\mathcal{B}(X)$, the first definition of regularity is automatic. These are always outer and inner regular provided inner regularity refers to closed sets. Note that if $A \supseteq B$ then $A \setminus B = B^C \setminus A^C$.

Lemma 9.8.4 Let μ be a finite measure defined on a σ algebra $\mathcal{F} \supseteq \mathcal{B}(X)$ where X is a metric space. Then the following hold.

1. μ is regular on $\mathcal{B}(X)$ meaning 9.18, 9.19 whenever $F \in \mathcal{B}(X)$.
2. μ is outer regular satisfying 9.19 on sets of \mathcal{F} if and only if it is inner regular satisfying 9.18 on sets of \mathcal{F} .
3. If μ is either inner or outer regular on sets of \mathcal{F} then if E is any set of \mathcal{F} , there exist F an F_σ set and G a G_δ set such that $F \subseteq E \subseteq G$ and $\mu(G \setminus F) = 0$.

Proof: 1.) First note every open set is the countable union of closed sets and every closed set is the countable intersection of open sets. Here is why. Let V be an open set and let

$$K_k \equiv \{x \in V : \text{dist}(x, V^C) \geq 1/k\}.$$

Then clearly the union of the K_k equals V . Thus

$$\mu(V) = \sup \{\mu(K) : K \subseteq V \text{ and } K \text{ is closed}\}.$$

If U is open and contains V , then $\mu(U) \geq \mu(V)$ and so

$$\mu(V) \leq \inf \{\mu(U) : U \supseteq V, U \text{ open}\} \leq \mu(V) \text{ since } V \subseteq V.$$

Thus μ is inner and outer regular on open sets. In what follows, K will be closed and V will be open.

Let \mathcal{H} be the open sets. This is a π system since it is closed with respect to finite intersections. Let

$$\mathcal{G} \equiv \{E \in \mathcal{B}(X) : \mu \text{ is inner and outer regular on } E\} \text{ so } \mathcal{G} \supseteq \mathcal{H}.$$

For $E \in \mathcal{G}$, let $V \supseteq E \supseteq K$ such that $\mu(V \setminus K) = \mu(V \setminus E) + \mu(E \setminus K) < \varepsilon$. Thus $K^C \supseteq E^C$ and so $\mu(K^C \setminus E^C) = \mu(E \setminus K) < \varepsilon$. Thus μ is outer regular on E^C because

$$\mu(K^C) = \mu(E^C) + \mu(K^C \setminus E^C) < \mu(E^C) + \varepsilon, K^C \supseteq E^C$$

Also, $E^C \supseteq V^C$ and $\mu(E^C \setminus V^C) = \mu(V \setminus E) < \varepsilon$ so μ is inner regular on E^C and so \mathcal{G} is closed for complements. If the sets of \mathcal{G} $\{E_i\}$ are disjoint, let $K_i \subseteq E_i \subseteq V_i$ with $\mu(V_i \setminus K_i) < \varepsilon 2^{-i}$. Then for $E \equiv \cup_i E_i$, and choosing m sufficiently large,

$$\mu(E) = \sum_i \mu(E_i) \leq \sum_{i=1}^m \mu(E_i) + \varepsilon \leq \sum_{i=1}^m \mu(K_i) + 2\varepsilon = \mu(\cup_{i=1}^m K_i) + 2\varepsilon$$

and so μ is inner regular on $E \equiv \cup_i E_i$. It remains to show that μ is outer regular on E . Letting $V \equiv \cup_i V_i$,

$$\mu(V \setminus E) \leq \mu(\cup_i (V_i \setminus E_i)) \leq \sum_i \varepsilon 2^{-i} = \varepsilon.$$

Hence μ is outer regular on E since $\mu(V) = \mu(E) + \mu(V \setminus E) \leq \mu(E) + \varepsilon$ and $V \supseteq E$.

By Dynkin's lemma, $\mathcal{G} = \sigma(\mathcal{H}) \equiv \mathcal{B}(X)$.

2.) Suppose that μ is outer regular on sets of $\mathcal{F} \supseteq \mathcal{B}(X)$. Letting $E \in \mathcal{F}$, by outer regularity, there exists an open set $V \supseteq E^C$ such that $\mu(V) - \mu(E^C) < \varepsilon$. Since μ is finite, $\varepsilon > \mu(V) - \mu(E^C) = \mu(V \setminus E^C) = \mu(E \setminus V^C) = \mu(E) - \mu(V^C)$ and V^C is a closed set contained in E . Therefore, if 9.19 holds, then so does 9.18. The converse is proved in the same way.

3.) The last claim is obtained by letting $G = \cap_n V_n$ where V_n is open, contains E , $V_n \supseteq V_{n+1}$, and $\mu(V_n) < \mu(E) + \frac{1}{n}$ and K_n , increasing closed sets contained in E such that $\mu(E) < \mu(K_n) + \frac{1}{n}$. Then let $F \equiv \cup_n F_n$ and $G \equiv \cap_n V_n$. Then $F \subseteq E \subseteq G$ and $\mu(G \setminus F) \leq \mu(V_n \setminus K_n) < 2/n$. ■

Next is a lemma which allows the replacement of closed with compact in the definition of inner regular.

Lemma 9.8.5 *Let μ be a finite measure on a σ algebra containing $\mathcal{B}(X)$, the Borel sets of X , a separable complete metric space, Polish space. Then if C is a closed set,*

$$\mu(C) = \sup \{ \mu(K) : K \subseteq C \text{ and } K \text{ is compact.} \}$$

It follows that for a finite measure on $\mathcal{B}(X)$ where X is a Polish space, μ is inner regular in the sense that for all $F \in \mathcal{B}(X)$,

$$\mu(F) = \sup \{ \mu(K) : K \subseteq F \text{ and } K \text{ is compact} \}$$

Proof: Let $\{a_k\}$ be a countable dense subset of C . Thus $\bigcup_{k=1}^{\infty} B(a_k, \frac{1}{n}) \supseteq C$. Therefore, there exists m_n such that

$$\mu \left(C \setminus \overline{\bigcup_{k=1}^{m_n} B \left(a_k, \frac{1}{n} \right)} \right) \equiv \mu(C \setminus C_n) < \frac{\varepsilon}{2^n}, \quad \overline{\bigcup_{k=1}^{m_n} B \left(a_k, \frac{1}{n} \right)} \equiv C_n.$$

Now let $K = C \cap (\bigcap_{n=1}^{\infty} C_n)$. Then K is a subset of C_n for each n and so for each $\varepsilon > 0$ there exists an ε net for K since C_n has a $1/n$ net, namely a_1, \dots, a_{m_n} . Since K is closed, it is complete and so it is also compact since it is complete and totally bounded, Theorem 3.5.8. Now

$$\mu(C \setminus K) \leq \mu \left(\bigcup_{n=1}^{\infty} (C \setminus C_n) \right) < \sum_{n=1}^{\infty} \frac{\varepsilon}{2^n} = \varepsilon.$$

Thus $\mu(C)$ can be approximated by $\mu(K)$ for K a compact subset of C . The last claim follows from Lemma 9.8.4. ■

The next theorem is the main result. It says that if the measure is outer regular and μ is σ finite then there is an approximation for $E \in \mathcal{F}$ in terms of F_{σ} and G_{δ} sets in which the F_{σ} set is a countable union of compact sets. Also μ is inner and outer regular on \mathcal{F} .

Theorem 9.8.6 *Suppose (X, \mathcal{F}, μ) , $\mathcal{F} \supseteq \mathcal{B}(X)$ is a measure space for X a metric space and μ is σ finite, $X = \bigcup_n X_n$ with $\mu(X_n) < \infty$ and the X_n disjoint. Suppose also that μ is outer regular. Then for each $E \in \mathcal{F}$, there exists F, G an F_{σ} and G_{δ} set respectively such that $F \subseteq E \subseteq G$ and $\mu(G \setminus F) = 0$. In particular, μ is inner and outer regular on \mathcal{F} . In case X is a complete separable metric space (Polish space), one can have F in the above be the countable union of compact sets and μ is inner regular in the sense of 9.20.*

Proof: Since μ is outer regular and $\mu(X_n) < \infty$, there exists an open set $V_n \supseteq E \cap X_n$ such that

$$\mu(V_n \setminus (E \cap X_n)) = \mu(V_n) - \mu(E \cap X_n) < \frac{\varepsilon}{2^n}.$$

Then let $V \equiv \bigcup_n V_n$ so that $V \supseteq E$. Then $E = \bigcup_n E \cap X_n$ and so

$$\mu(V \setminus E) \leq \mu \left(\bigcup_n (V_n \setminus (E \cap X_n)) \right) \leq \sum_n \mu(V_n \setminus (E \cap X_n)) < \sum_n \frac{\varepsilon}{2^n} = \varepsilon$$

Similarly, there exists U_n open such that $\mu(U_n \setminus (E^C \cap X_n)) < \frac{\varepsilon}{2^n}$, $U_n \supseteq E^C \cap X_n$ so if $U \equiv \bigcup_n U_n$, $\mu(U \setminus E^C) = \mu(E \setminus U^C) < \varepsilon$. Now U^C is closed and contained in E because $U \supseteq E^C$. Hence, letting $\varepsilon = \frac{1}{2^n}$, there exist closed sets C_n , and open sets V_n such that $C_n \subseteq E \subseteq V_n$ and $\mu(V_n \setminus C_n) < \frac{1}{2^{n-1}}$. Letting $G \equiv \bigcap_n V_n$, $F \equiv \bigcup_n C_n$, $F \subseteq E \subseteq G$ and $\mu(G \setminus F) \leq \mu(V_n \setminus C_n) < \frac{1}{2^{n-1}}$. Since n is arbitrary, $\mu(G \setminus F) = 0$.

To finish the proof, I will use Lemma 9.8.5 in the case where X is a Polish space.

By the first part, $\mu(G \setminus F) = 0$ where F is the countable union of closed sets $\{C_k\}_{k=1}^\infty$ and $F \subseteq E \subseteq G$. Letting $\mu_n(E) \equiv \mu(E \cap X_n)$, μ_n is a finite measure and so if C_k is one of those closed sets, Lemma 9.8.5 implies

$$\mu_n(C_k) \equiv \mu(C_k \cap X_n) = \sup \{ \mu(K \cap X_n) : K \subseteq C_k, K \text{ compact} \}$$

Pick K_k compact such that $\mu_n(C_k \setminus K_k) < \frac{\varepsilon}{2^k}$, $K_k \subseteq C_k$. Then letting $\hat{F} \equiv \cup_k K_k$, it follows \hat{F} is a countable union of compact sets contained in F and

$$\mu(F \setminus \hat{F}) = \mu(\cup_k C_k \setminus (\cup_k K_k)) \leq \mu(\cup_k (C_k \setminus K_k)) \leq \sum_k \mu(C_k \setminus K_k) < \varepsilon$$

Therefore, letting \hat{F}_m be a countable union of compact sets contained in F for which $\mu(F \setminus \hat{F}_m) < \frac{1}{2^m}$, let $\tilde{F} \equiv \cup_m \hat{F}_m$. Then \tilde{F} is a countable union of compact sets and

$$\mu(F \setminus \tilde{F}) \leq \mu(F \setminus \hat{F}_m) < \frac{1}{2^m}$$

and so $\mu(F \setminus \tilde{F}) = 0$. Then

$$\mu(G \setminus \tilde{F}) = \mu(G \setminus F) + \mu(F \setminus \tilde{F}) = \mu(G \setminus F) = 0$$

so as claimed, one can have F in the first part be the countable union of compact sets. Letting $E \in \mathcal{F}$, it was just shown that there exist G a G_δ set and F the countable union of compact sets such that $\mu(G \setminus F) = 0$, $F \subseteq E \subseteq G$. Therefore, $\mu(E) = \mu(E \setminus F) + \mu(F) = \mu(F)$ and so this shows inner regularity in the sense of 9.20 because if $l < \mu(E) = \mu(F)$, one could include enough of the compact sets whose union is F to obtain a compact set K for which $\mu(K) > l$. ■

An important example is the case of a random vector and its distribution measure.

Definition 9.8.7 A measurable function $X : (\Omega, \mathcal{F}, \mu) \rightarrow Z$ a metric space is called a random variable when $\mu(\Omega) = 1$. For such a random variable, one can define a distribution measure λ_X on the Borel sets of Z as follows. $\lambda_X(G) \equiv \mu(X^{-1}(G))$. This is a well defined measure on the Borel sets of Z because it makes sense for every G open and $\mathcal{G} \equiv \{G \subseteq Z : X^{-1}(G) \in \mathcal{F}\}$ is a σ algebra which contains the open sets, hence the Borel sets. Such a random variable is also called a random vector when Z is a vector space.

Corollary 9.8.8 Let X be a random variable with values in a separable complete metric space Z . Then λ_X is an inner and outer regular measure defined on $\mathcal{B}(Z)$.

One such example of a complete metric space and a measure which is finite on compact sets is the following where the closures of balls are compact. Thus, this involves finite dimensional situations essentially. Note that if you have a metric space in which the closures of balls are compact sets, then the metric space must be separable. This is because you can pick a point ξ and consider the closures of balls $\overline{B(\xi, n)}$. Then $\overline{B(\xi, n)}$ is complete and totally bounded so it has a countable dense subset D_n . Let $D = \cup_n D_n$.

Corollary 9.8.9 Let Ω be a complete metric space which is the countable union of compact sets K_n and suppose, for μ a Borel measure, $\mu(K_n)$ is finite. Then μ must be regular on $\mathcal{B}(\Omega)$. In particular, if Ω is a metric space and the closure of each ball is compact, and μ is finite on balls, then μ must be regular.

Proof: Let the compact sets be increasing without loss of generality, and let $\mu_n(E) \equiv \mu(K_n \cap E)$. Thus μ_n is a finite measure defined on the Borel sets of a Polish space so it is regular. Letting $l < \mu(E)$, there exists n such that $l < \mu_n(E) \leq \mu(E)$. By what was shown above in Lemma 9.8.5, there exists H compact, $H \subseteq E$ such that also for a large n , $\mu_n(H) > l$. Hence $\mu(H \cap K_n) > l$ and so μ is inner regular. It remains to verify that μ is outer regular. If $\mu(E) = \infty$, there is nothing to show. Assume then that $\mu(E) < \infty$. Let $V_n \supseteq E$ with $\mu_n(V_n \setminus E) < \varepsilon 2^{-n}$ so also $\mu(V_n) < \infty$. We can assume also that $V_n \supseteq V_{n+1}$ for all n . Thus $\mu((V_n \setminus E) \cap K_n) < 2^{-n}\varepsilon$. Let $G = \bigcap_k V_k$. Then $G \subseteq V_n$ so $\mu((G \setminus E) \cap K_n) < 2^{-n}\varepsilon$. Letting $n \rightarrow \infty$, $\mu(G \setminus E) = 0$ and $G \supseteq E$. Then, since V_1 has finite measure, $\mu(G \setminus E) = \lim_{n \rightarrow \infty} \mu(V_n \setminus E)$ and so for all n large enough, $\mu(V_n \setminus E) < \varepsilon$ so $\mu(E) + \varepsilon > \mu(V_n)$ and so μ is outer regular. In the last case, if the closure of each ball is compact, then Ω is automatically complete because every Cauchy sequence is contained in some ball and so has a convergent subsequence. Since the sequence is Cauchy, it also converges by Theorem 3.2.2 on Page 73. ■

9.9 One Dimensional Lebesgue Stieltjes Measure

Now with these major results about measures, it is time to specialize to the outer measure of Theorem 9.7.4. The next theorem gives Lebesgue Stieltjes measure on \mathbb{R} . The conditions 9.21 and 9.22 given below are known respectively as inner and outer regularity.

Theorem 9.9.1 *Let \mathcal{F} denote the σ algebra of Theorem 9.5.4, associated with the outer measure μ in Theorem 9.7.4, on which μ is a measure. Then every open interval is in \mathcal{F} . So are all open and closed sets and consequently all Borel sets. Furthermore, if E is any set in \mathcal{F}*

$$\mu(E) = \sup\{\mu(K) : K \text{ compact, } K \subseteq E\} \quad (9.21)$$

$$\mu(E) = \inf\{\mu(V) : V \text{ is an open set } V \supseteq E\} \quad (9.22)$$

If $E \in \mathcal{F}$, there exists F a countable union of compact sets, an F_σ set and a set G a countable intersection of open sets, a G_δ set such that $F \subseteq E \subseteq G$ but $\mu(G \setminus F) = 0$. Also μ is finite on compact sets.

Proof: By Theorem 9.7.4 and Theorem 9.6.1 the σ algebra includes the Borel sets $\mathcal{B}(\mathbb{R})$. However, note that \mathcal{F} is complete and there is no such requirement for this measure on $\mathcal{B}(\mathbb{R})$. Thus it is reasonable to think that \mathcal{F} could be larger than $\mathcal{B}(\mathbb{R})$.

Now consider the last claim about regularity. The assertion of outer regularity on \mathcal{F} is not hard to get. Letting E be any set $\mu(E) < \infty$, there exist open intervals covering E denoted by $\{(a_i, b_i)\}_{i=1}^\infty$ such that

$$\mu(E) + \varepsilon > \sum_{i=1}^\infty F(b_i -) - F(a_i +) = \sum_{i=1}^\infty \mu(a_i, b_i) \geq \mu(V)$$

where V is the union of the open intervals just mentioned. Thus

$$\mu(E) \leq \mu(V) \leq \mu(E) + \varepsilon.$$

This shows outer regularity. If $\mu(E) = \infty$, there is nothing to show. Since μ is finite on intervals, it is σ finite. It follows from Theorem 9.8.6 that μ is inner regular also and the claim about approximation with F_σ and G_δ sets follows. ■

Definition 9.9.2 When the integrator function is $F(x) = x$, the Lebesgue Stieltjes measure just discussed is known as one dimensional Lebesgue measure and is denoted as m .

Proposition 9.9.3 For m Lebesgue measure, $m([a, b]) = m((a, b)) = b - a$. Also m is translation invariant in the sense that if E is any Lebesgue measurable set, then $m(x + E) = m(E)$.

Proof: The formula for the measure of an interval comes right away from Theorem 9.7.4. From this, it follows right away that whenever E is an interval, $m(x + E) = m(E)$. Every open set is the countable disjoint union of open intervals, so if E is an open set, then $m(x + E) = m(E)$. What about closed sets? First suppose H is a closed and bounded set. Then letting $(-n, n) \supseteq H$,

$$\mu(((-n, n) \setminus H) + x) + \mu(H + x) = \mu((-n, n) + x)$$

Hence, from what was just shown about open sets,

$$\begin{aligned} \mu(H) &= \mu((-n, n)) - \mu((-n, n) \setminus H) \\ &= \mu((-n, n) + x) - \mu(((-n, n) \setminus H) + x) = \mu(H + x) \end{aligned}$$

Therefore, the translation invariance holds for closed and bounded sets. If H is an arbitrary closed set, then

$$\mu(H + x) = \lim_{n \rightarrow \infty} \mu(H \cap [-n, n] + x) = \lim_{n \rightarrow \infty} \mu(H \cap [-n, n]) = \mu(H).$$

It follows right away that μ is translation invariant on F_σ and G_δ sets. Now using Theorem 9.9.1, if E is an arbitrary measurable set, there exist an F_σ set F and a G_δ set G such that $F \subseteq E \subseteq G$ and $m(F) = m(G) = m(E)$. Then

$$m(F) = m(x + F) \leq m(x + E) \leq m(x + G) = m(G) = m(E) = m(F). \quad \blacksquare$$

9.10 Exercises

1. Show carefully that if \mathfrak{S} is a set whose elements are σ algebras which are subsets of $\mathcal{P}(\Omega)$, then $\cap \mathfrak{S}$ is also a σ algebra. Now let $\mathcal{G} \subseteq \mathcal{P}(\Omega)$ satisfy property P if \mathcal{G} is closed with respect to complements and countable disjoint unions as in Dynkin's lemma, and contains \emptyset and Ω . If $\mathfrak{H} \subseteq \mathcal{G}$ is any set whose elements are subsets of $\mathcal{P}(\Omega)$ which satisfies property P , then $\cap \mathfrak{H}$ also satisfies property P . Thus there is a smallest subset of \mathcal{G} satisfying P . In other words, verify the details of the proof of Dynkin's lemma.
2. The Borel sets of a metric space (X, d) are the sets in the smallest σ algebra which contains the open sets. These sets are denoted as $\mathcal{B}(X)$. Thus $\mathcal{B}(X) = \sigma(\text{open sets})$ where $\sigma(\mathcal{F})$ simply means the smallest σ algebra which contains \mathcal{F} . Show that in \mathbb{R}^n , $\mathcal{B}(\mathbb{R}^n) = \sigma(\mathcal{P})$ where \mathcal{P} consists of the half open rectangles which are of the form $\prod_{i=1}^n [a_i, b_i)$.
3. Recall that $f : (\Omega, \mathcal{F}) \rightarrow X$ where X is a metric space is measurable means that inverse images of open sets are in \mathcal{F} . Show that if E is any set in $\mathcal{B}(X)$, then

$f^{-1}(E) \in \mathcal{F}$. Thus, inverse images of Borel sets are measurable. Next consider $f : (\Omega, \mathcal{F}) \rightarrow X$ being measurable and $g : X \rightarrow Y$ is Borel measurable, meaning that $g^{-1}(\text{open}) \in \mathcal{B}(X)$. Explain why $g \circ f$ is measurable. **Hint:** You know that $(g \circ f)^{-1}(U) = f^{-1}(g^{-1}(U))$. For your information, it does not work the other way around. That is, measurable composed with Borel measurable is not necessarily measurable. In fact examples exist which show that if g is measurable and f is continuous, then $g \circ f$ may fail to be measurable. An example is given later.

4. If you have X_i is a metric space, let $X = \prod_{i=1}^n X_i$ with the metric

$$d(x, y) \equiv \max\{d_i(x_i, y_i), i = 1, 2, \dots, n\}$$

You considered this in an earlier problem. Show that any set of the form

$$\prod_{i=1}^n E_i, E_i \in \mathcal{B}(X_i)$$

is a Borel set. That is, the product of Borel sets is Borel. **Hint:** You might consider the continuous functions $\pi_i : \prod_{j=1}^n X_j \rightarrow X_i$ which are the projection maps. Thus $\pi_i(x) \equiv x_i$. Then $\pi_i^{-1}(E_i)$ would have to be Borel measurable whenever $E_i \in \mathcal{B}(X_i)$. Explain why. You know π_i is continuous. Why would $\pi_i^{-1}(\text{Borel})$ be a Borel set? Then you might argue that $\prod_{i=1}^n E_i = \cap_{i=1}^n \pi_i^{-1}(E_i)$.

5. You have two finite measures defined on $\mathcal{B}(X)$ μ, ν . Suppose these are equal on every open set. Show that these must be equal on every Borel set. **Hint:** You should use Dynkin's lemma to show this very easily.
6. Show that $(\mathbb{N}, \mathcal{P}(\mathbb{N}), \mu)$ is a measure space where $\mu(S)$ equals the number of elements of S . You need to verify that if the sets E_i are disjoint, then $\mu(\cup_{i=1}^{\infty} E_i) = \sum_{i=1}^{\infty} \mu(E_i)$.
7. Let Ω be an uncountable set and let \mathcal{F} denote those subsets of Ω , F such that either F or F^C is countable. Show that this is a σ algebra. Next define the following measure. $\mu(A) = 1$ if A is uncountable and $\mu(A) = 0$ if A is countable. Show that μ is a measure. This is a perverted example.
8. Let $\mu(E) = 1$ if $0 \in E$ and $\mu(E) = 0$ if $0 \notin E$. Show this is a measure on $\mathcal{P}(\mathbb{R})$.
9. Give an example of a measure μ and a measure space and a decreasing sequence of measurable sets $\{E_i\}$ such that $\lim_{n \rightarrow \infty} \mu(E_n) \neq \mu(\cap_{i=1}^{\infty} E_i)$.
10. You have a measure space (Ω, \mathcal{F}, P) where P is a probability measure on \mathcal{F} . Then you also have a measurable function $X : \Omega \rightarrow Z$ where Z is some metric space. Thus $X^{-1}(U) \in \mathcal{F}$ whenever U is open. Now define a measure on $\mathcal{B}(Z)$ denoted by λ_X and defined by $\lambda_X(E) = P(\{\omega : X(\omega) \in E\})$. Explain why this yields a well defined probability measure on $\mathcal{B}(Z)$ which is regular. This is called the distribution measure.
11. Let $K \subseteq V$ where K is closed and V is open. Consider the following function.

$$f(x) = \frac{\text{dist}(x, V^C)}{\text{dist}(x, K) + \text{dist}(x, V^C)}$$

Explain why this function is continuous, equals 0 off V and equals 1 on K .

12. Let (Ω, \mathcal{F}) be a measurable space and let $f : \Omega \rightarrow X$ be a measurable function. Then $\sigma(f)$ denotes the smallest σ algebra such that f is measurable with respect to this σ algebra. Show that $\sigma(f) = \{f^{-1}(E) : E \in \mathcal{B}(X)\}$.
13. Let $(\Omega, \mathcal{F}, \mu)$ be a measure space. A sequence of functions $\{f_n\}$ is said to converge in measure to a measurable function f if and only if for each

$$\varepsilon > 0, \lim_{n \rightarrow \infty} \mu(\omega : |f_n(\omega) - f(\omega)| > \varepsilon) = 0.$$

Show that if this happens, then there exists a subsequence $\{f_{n_k}\}$ and a set of measure N such that if $\omega \notin N$, then $\lim_{k \rightarrow \infty} f_{n_k}(\omega) = f(\omega)$. Also show that if $\lim_{n \rightarrow \infty} f_n(\omega) = f(\omega)$, and $\mu(\Omega) < \infty$, then f_n converges in measure to f . **Hint:** For the subsequence, let $\mu(\omega : |f_{n_k}(\omega) - f(\omega)| > \varepsilon) < 2^{-k}$ and use Borel Cantelli lemma.

14. Let X, Y be separable metric spaces. Then $X \times Y$ can also be considered as a metric space with the metric $\rho((x, y), (\hat{x}, \hat{y})) \equiv \max(d_X(x, \hat{x}), d_Y(y, \hat{y}))$. Verify this. Then show that if \mathcal{H} consists of sets $A \times B$ where A, B are Borel sets in X and Y respectively, then it follows that $\sigma(\mathcal{H}) = \mathcal{B}(X \times Y)$, the Borel sets from $X \times Y$. Extend to the Cartesian product $\prod_i X_i$ of finitely many separable metric spaces.

9.11 Completion of a Measure Space

Next is the notion of the completion of a measure space. The idea is that you might not have completeness in your measure space but you can always complete it.

Definition 9.11.1 Recall that a measure space $(\Omega, \mathcal{F}, \lambda)$ is σ finite if there is a countable set $\{\Omega_n\}_{n=1}^{\infty}$ such that $\cup_n \Omega_n = \Omega$ and $\lambda(\Omega_n) < \infty$.

The next theorem is like some earlier ones related to regularity including the approximation with G_δ and F_σ sets. The arguments are similar.

Theorem 9.11.2 Let $(\Omega, \mathcal{F}, \mu)$ be a measure space. Then there exists a measure space, $(\Omega, \mathcal{G}, \lambda)$ satisfying

1. $(\Omega, \mathcal{G}, \lambda)$ is a complete measure space.
2. $\lambda = \mu$ on \mathcal{F}
3. $\mathcal{G} \supseteq \mathcal{F}$
4. For every $E \in \mathcal{G}$ there exists $G \in \mathcal{F}$ such that $G \supseteq E$ and $\mu(G) = \lambda(E)$.

In addition to this, if $(\Omega, \mathcal{F}, \mu)$ is σ finite, then the following approximation result holds.

5. For every $E \in \mathcal{G}$ there exists $F \in \mathcal{F}$ and $G \in \mathcal{F}$ such that $F \subseteq E \subseteq G$ and

$$\mu(G \setminus F) = \lambda(G \setminus F) = 0 \tag{9.23}$$

There is a unique complete measure space $(\Omega, \mathcal{G}, \lambda)$ extending $(\Omega, \mathcal{F}, \mu)$ which satisfies 9.23. In particular, there are no new sets if the original measure space was already complete.

Proof: Define the outer measure

$$\lambda(A) \equiv \inf \{ \mu(E) : E \in \mathcal{F}, E \supseteq A \}, \lambda(\emptyset) \equiv 0.$$

Denote by \mathcal{G} the σ algebra of λ measurable sets. Then $(\Omega, \mathcal{G}, \lambda)$ is complete by the general Caratheodory procedure presented earlier.

I claim that $\lambda = \mu$ on \mathcal{F} . If $A \in \mathcal{F}$,

$$\mu(A) \leq \inf \{ \mu(E) : E \in \mathcal{F}, E \supseteq A \} \equiv \lambda(A) \leq \mu(A)$$

because $A \supseteq A$. Thus, these are all equal in the above and $\lambda = \mu$ on \mathcal{F} .

Why is $\mathcal{F} \subseteq \mathcal{G}$? Letting $\lambda(S) < \infty$, (There is nothing to prove if $\lambda(S) = \infty$.) let $G \in \mathcal{F}$ be such that $G \supseteq S$ and $\lambda(S) = \mu(G)$. This is possible because

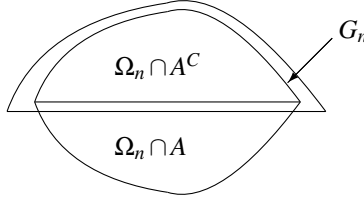
$$\lambda(S) \equiv \inf \{ \mu(E) : E \supseteq S \text{ and } E \in \mathcal{F} \}.$$

Then if $A \in \mathcal{F}$,

$$\begin{aligned} \lambda(S) &\leq \lambda(S \cap A) + \lambda(S \cap A^C) \leq \lambda(G \cap A) + \lambda(G \cap A^C) \\ &= \mu(G \cap A) + \mu(G \cap A^C) = \mu(G) = \lambda(S). \end{aligned}$$

Thus $\mathcal{F} \subseteq \mathcal{G}$.

Finally suppose μ is σ finite. Let $\Omega = \bigcup_{n=1}^{\infty} \Omega_n$ where the Ω_n are disjoint sets of \mathcal{F} and $\mu(\Omega_n) < \infty$. If the Ω_n are not disjoint, replace Ω_n with $\Omega_n \setminus \bigcup_{k=1}^{n-1} \Omega_k$. Letting $A \in \mathcal{G}$, consider $A_n \equiv A \cap \Omega_n$. From what was just shown, there exists $G_n \supseteq A^C \cap \Omega_n$, $G_n \subseteq \Omega_n$ such that $\mu(G_n) = \lambda(A^C \cap \Omega_n)$, $G_n \in \mathcal{F}$.



Since $\mu(\Omega_n) < \infty$, this implies

$$\lambda(G_n \setminus (A^C \cap \Omega_n)) = \lambda(G_n) - \lambda(A^C \cap \Omega_n) = 0.$$

Now $G_n^C \subseteq A \cup \Omega_n^C$ but $G_n \subseteq \Omega_n$ and so $G_n^C \subseteq A \cup \Omega_n$. Define $F_n \equiv G_n^C \cap \Omega_n \subseteq A_n$ and it follows $\lambda(A_n \setminus F_n) =$

$$\begin{aligned} \lambda(A \cap \Omega_n \setminus (G_n^C \cap \Omega_n)) &= \lambda(A \cap \Omega_n \cap G_n) = \lambda(A \cap G_n) = \lambda(G_n \setminus A^C) \\ &\leq \lambda(G_n \setminus (A^C \cap \Omega_n)) = 0. \end{aligned}$$

Letting $F = \bigcup_n F_n$, it follows that $F \in \mathcal{F}$ and

$$\lambda(A \setminus F) \leq \sum_{k=1}^{\infty} \lambda(A_k \setminus F_k) = 0.$$

Also, there exists $G_n \supseteq A_n$ such that $\mu(G_n) = \lambda(G_n) = \lambda(A_n)$. Since the measures are finite, it follows that $\lambda(G_n \setminus A_n) = 0$. Then letting $G = \bigcup_{n=1}^{\infty} G_n$, it follows that $G \supseteq A$ and

$$\begin{aligned} \lambda(G \setminus A) &= \lambda(\bigcup_{n=1}^{\infty} G_n \setminus \bigcup_{n=1}^{\infty} A_n) \\ &\leq \lambda(\bigcup_{n=1}^{\infty} (G_n \setminus A_n)) \leq \sum_{n=1}^{\infty} \lambda(G_n \setminus A_n) = 0. \end{aligned}$$

Thus $\mu(G \setminus F) = \lambda(G \setminus F) = \lambda(G \setminus A) + \lambda(A \setminus F) = 0$.

If you have (λ', \mathcal{G}') complete and satisfying 9.23, then letting $E \in \mathcal{G}'$, it follows from 5, that there exist $F, G \in \mathcal{F}$ such that

$$F \subseteq E \subseteq G, \mu(G \setminus F) = 0 = \lambda(G \setminus F).$$

Therefore, by completeness of the two measure spaces, $E \in \mathcal{G}$. The opposite inclusion is similar. Hence $\mathcal{G} = \mathcal{G}'$. If $E \in \mathcal{G}$, let $F \subseteq E \subseteq G$ where $\mu(G \setminus F) = 0$. Then

$$\lambda(E) \leq \mu(G) = \mu(F) = \lambda'(F) \leq \lambda'(E)$$

The opposite inequality holds by the same reasoning. Hence $\lambda = \lambda'$. If $(\Omega, \mathcal{F}, \mu)$ is already complete, then you could let this be $(\Omega, \mathcal{G}', \lambda')$ and find that $\mathcal{G} = \mathcal{F} = \mathcal{G}'$. ■

Another useful result is the following.

Corollary 9.11.3 *Suppose, in the situation of Theorem 9.11.2, $f \geq 0$ and is \mathcal{G} measurable. Then there exists $g, 0 \leq g \leq f$ and $f = g$ for all ω off a set of measure zero.*

Proof: Let $s_n \uparrow f$ where $s_n(\omega) = \sum_{i=1}^{m_n} c_i \mathcal{X}_{E_i}(\omega)$ for $\omega \in \Omega$. Then by the regularity assertion of this theorem, there exists $F_i \in \mathcal{F}$ such that $F_i \subseteq E_i$ and $\lambda(E_i \setminus F_i) = 0$. Then let $\hat{s}_n(\omega) = \sum_{i=1}^{m_n} c_i \mathcal{X}_{F_i}(\omega)$. Then $\hat{s}_n \leq s_n$ and letting $N = \bigcup_{n=1}^{\infty} \{\omega : s_n(\omega) \neq \hat{s}_n(\omega)\}$, it follows that $\lambda(N) = 0$ and for $\omega \notin N$,

$$\hat{s}_n(\omega) = s_n(\omega) \rightarrow f(\omega) = g(\omega).$$

Now let $g(\omega) \equiv \liminf_{n \rightarrow \infty} \hat{s}_n(\omega) \leq \lim_{n \rightarrow \infty} s_n(\omega) = f(\omega)$ and g is \mathcal{F} measurable because if $g_n(\omega) = \inf\{\hat{s}_k : k \geq n\}$, this is \mathcal{F} measurable since

$$g_n^{-1}((-\infty, a)) = \bigcup_{k \geq n} \hat{s}_k^{-1}((-\infty, a)) \in \mathcal{F}$$

Now g being the limit of these g_n , it follows that g is also \mathcal{F} measurable. ■

This will show that in most situations, you can simply modify your function on a set of measure zero and consider one which is \mathcal{F} measurable.

Recall Corollary 9.8.9 about regularity. Then there is an easy corollary.

Corollary 9.11.4 *Let Ω be a complete metric space which is the countable union of compact sets K_n and suppose, for μ a Borel measure, $\mu(K_n)$ is finite. Then μ must be regular on $\mathcal{B}(\Omega)$. If $(\bar{\mu}, \mathcal{G})$ is the completion, then $\bar{\mu}$ is inner and outer regular on sets of \mathcal{G} . Also, if $E \in \mathcal{G}$, there are F_σ and G_δ sets, F, G respectively such that $\bar{\mu}(G \setminus F) = 0$ and $F \subseteq E \subseteq G$.*

9.12 Vitali Coverings

There is another covering theorem which may also be referred to as the Besicovitch covering theorem. At first, the balls will be closed but this assumption will be removed. Assume the following: $(X, \bar{\mu})$ is a finite dimensional normed linear space of dimension p and $\bar{\mu}$ is an outer measure on $\mathcal{P}(X)$. We really have in mind that X is \mathbb{R}^p with some norm. Assume the following:

1. Let μ the measure determined by $\bar{\mu}$ on the σ algebra \mathcal{S} which contains the Borel sets.

2. Or let μ be a measure on \mathcal{S} where \mathcal{S} contains the Borel sets and $\bar{\mu}$ is the outer measure determined by μ as described in Proposition 9.4.2. Always assume the following:
3. $\mu(B(x, r)) < \infty$.
4. If $E \in \mathcal{S}$, then

$$\begin{aligned}\mu(E) &= \sup\{\mu(K) : K \subseteq E \text{ and } K \text{ is compact}\} \\ \mu(E) &= \inf\{\mu(V) : V \supseteq E \text{ and } V \text{ is open}\}\end{aligned}$$

If this measure μ is also complete, then recall it is termed a Radon measure.

Note that \mathcal{S} is given to contain all closed sets and open sets. The above situation is very common. See Corollary 9.8.9 which gives 4 follows from 3. In fact, the above is the typical case for measures on finite dimensional spaces.

Definition 9.12.1 A collection of balls, \mathcal{F} covers a set E in the sense of Vitali if whenever $x \in E$ and $\varepsilon > 0$, there exists a ball $B \in \mathcal{F}$ whose center is x having diameter less than ε .

I will give a proof of the following theorem.

Theorem 9.12.2 Let E be a set with $\bar{\mu}(E) < \infty$ and either 1 or 2 along with the regularity conditions 3 and 4. Suppose \mathcal{F} is a collection of closed balls which cover E in the sense of Vitali. Then there exists a sequence of disjoint balls $\{B_i\} \subseteq \mathcal{F}$ such that $\bar{\mu}(E \setminus \bigcup_{j=1}^N B_j) = 0$, $N \leq \infty$.

Proof: Let N_p be the constant of the Besicovitch covering theorem, Theorem 4.5.8. Choose $r > 0$ such that $(1-r)^{-1} \left(1 - \frac{1}{2N_p+2}\right) \equiv \lambda < 1$. If $\bar{\mu}(E) = 0$, there is nothing to prove so assume $\bar{\mu}(E) > 0$. Let U_1 be an open set containing E with $(1-r)\mu(U_1) < \bar{\mu}(E)$ and $2\bar{\mu}(E) > \mu(U_1)$, and let \mathcal{F}_1 be those sets of \mathcal{F} which are contained in U_1 whose centers are in E . Thus \mathcal{F}_1 is also a Vitali cover of E . Now by the Besicovitch covering theorem proved earlier, Theorem 4.5.8, there exist balls B , of \mathcal{F}_1 such that $E \subseteq \bigcup_{i=1}^{N_p} \{B : B \in \mathcal{G}_i\}$ where \mathcal{G}_i consists of a collection of disjoint balls of \mathcal{F}_1 . Therefore, $\bar{\mu}(E) \leq \sum_{i=1}^{N_p} \sum_{B \in \mathcal{G}_i} \mu(B)$ and so, for some $i \leq N_p$,

$$(N_p + 1) \sum_{B \in \mathcal{G}_i} \mu(B) > \bar{\mu}(E).$$

It follows there exists a finite set of balls of \mathcal{G}_i , $\{B_1, \dots, B_{m_1}\}$ such that

$$(N_p + 1) \sum_{i=1}^{m_1} \mu(B_i) > \bar{\mu}(E) \tag{9.24}$$

and so

$$(2N_p + 2) \sum_{i=1}^{m_1} \mu(B_i) > 2\bar{\mu}(E) > \mu(U_1).$$

Now 9.24 implies

$$\frac{\mu(U_1)}{2N_2+2} \leq \frac{2\bar{\mu}(E)}{2N_2+2} = \frac{\bar{\mu}(E)}{N_2+1} < \sum_{i=1}^{m_1} \mu(B_i).$$

Also U_1 was chosen such that $(1-r)\mu(U_1) < \bar{\mu}(E)$, and so

$$\begin{aligned} \lambda \bar{\mu}(E) &\geq \lambda(1-r)\mu(U_1) = \left(1 - \frac{1}{2N_p+2}\right) \mu(U_1) \\ &\geq \mu(U_1) - \sum_{i=1}^{m_1} \mu(B_i) = \mu(U_1) - \mu\left(\bigcup_{j=1}^{m_1} B_j\right) \\ &= \mu\left(U_1 \setminus \bigcup_{j=1}^{m_1} B_j\right) \geq \bar{\mu}\left(E \setminus \bigcup_{j=1}^{m_1} B_j\right). \end{aligned}$$

Since the balls are closed, you can consider the sets of \mathcal{F} which have empty intersection with $\bigcup_{j=1}^{m_1} B_j$ and this new collection of sets will be a Vitali cover of $E \setminus \bigcup_{j=1}^{m_1} B_j$. Letting this collection of balls play the role of \mathcal{F} in the above argument, and letting $E \setminus \bigcup_{j=1}^{m_1} B_j$ play the role of E , repeat the above argument and obtain disjoint sets of \mathcal{F} , $\{B_{m_1+1}, \dots, B_{m_2}\}$, such that

$$\lambda \bar{\mu}\left(E \setminus \bigcup_{j=1}^{m_1} B_j\right) > \bar{\mu}\left(\left(E \setminus \bigcup_{j=1}^{m_1} B_j\right) \setminus \bigcup_{j=m_1+1}^{m_2} B_j\right) = \bar{\mu}\left(E \setminus \bigcup_{j=1}^{m_2} B_j\right),$$

and so $\lambda^2 \bar{\mu}(E) > \bar{\mu}\left(E \setminus \bigcup_{j=1}^{m_2} B_j\right)$. Continuing in this way, yields a sequence of disjoint balls $\{B_i\}$ contained in \mathcal{F} and $\bar{\mu}\left(E \setminus \bigcup_{j=1}^N B_j\right) \leq \bar{\mu}\left(E \setminus \bigcup_{j=1}^{m_k} B_j\right) < \lambda^k \bar{\mu}(E)$ for all k . If the process stops because E gets covered, then N is finite and if not, then $N = \infty$. Therefore, $\bar{\mu}\left(E \setminus \bigcup_{j=1}^N B_j\right) = 0$ and this proves the Theorem. ■

It is not necessary to assume $\bar{\mu}(E) < \infty$. It is given that $\mu(B(x, R)) < \infty$. Letting $C(x, r)$ be all y with $\|y - x\| = r$. Then there are only finitely many $r < R$ such that $\mu(C(x, r)) \geq \frac{1}{n}$. Hence there are only countably many $r < R$ such that $\mu(C(x, r)) > 0$.

Corollary 9.12.3 *Let E nonempty set and either 1 or 2 along with the regularity conditions 3 and 4. Suppose \mathcal{F} is a collection of closed balls which cover E in the sense of Vitali. Then there exists a sequence of disjoint balls $\{B_i\} \subseteq \mathcal{F}$ such that*

$$\bar{\mu}\left(E \setminus \bigcup_{j=1}^N B_j\right) = 0, N \leq \infty$$

Proof: By 3, μ is finite on compact sets. Recall these are closed and bounded. There are at most countably many numbers, $\{b_i\}_{i=1}^{\infty}$ such that $\mu(C(0, b_i)) > 0$. It follows that there exists an increasing sequence of positive numbers, $\{r_i\}_{i=1}^{\infty}$ such that $\lim_{i \rightarrow \infty} r_i = \infty$ and $\mu(C(0, r_i)) = 0$. Now let

$$\begin{aligned} D_1 &\equiv \{x : \|x\| < r_1\}, D_2 \equiv \{x : r_1 < \|x\| < r_2\}, \\ \dots, D_m &\equiv \{x : r_{m-1} < \|x\| < r_m\}, \dots \end{aligned}$$

Let \mathcal{F}_m denote those closed balls of \mathcal{F} which are contained in D_m . Then letting E_m denote $E \cap D_m$, \mathcal{F}_m is a Vitali cover of E_m , $\bar{\mu}(E_m) < \infty$, and so by Theorem 9.12.2, there exists

a countable sequence of balls from $\mathcal{F}_m \left\{ B_j^m \right\}_{j=1}^N$, such that $\bar{\mu} \left(E_m \setminus \bigcup_{j=1}^N B_j^m \right) = 0$. Then consider the countable collection of balls, $\left\{ B_j^m \right\}_{j,m=1}^\infty$.

$$\begin{aligned} \bar{\mu} \left(E \setminus \bigcup_{m=1}^\infty \bigcup_{j=1}^N B_j^m \right) &\leq \bar{\mu} \left(\bigcup_{j=1}^\infty \partial B(\mathbf{0}, r_i) \right) + \\ + \sum_{m=1}^\infty \bar{\mu} \left(E_m \setminus \bigcup_{j=1}^N B_j^m \right) &= 0, N \leq \infty. \blacksquare \end{aligned}$$

If some E_m is empty, you could let your balls be the empty set.

You don't need to assume the balls are closed. In fact, the balls can be open, closed or anything in between and the same conclusion can be drawn provided you change the definition of a Vitali cover a little. For each point of the set covered, the covering includes all balls centered at that point having radius sufficiently small. In case that $\mu(C(\mathbf{x}, r)) = 0$ for all \mathbf{x}, r where $C(\mathbf{x}, r) \equiv \{\mathbf{y} : \|\mathbf{y} - \mathbf{x}\| = r\}$, no modification is necessary. This includes the case of Lebesgue measure. However, in the general case, consider the following modification of the notion of a Vitali cover.

Definition 9.12.4 Suppose \mathcal{F} is a collection of balls which cover E in the sense that for all $\varepsilon > 0$ there are uncountably many balls of \mathcal{F} centered at \mathbf{x} having radius less than ε .

Corollary 9.12.5 Let 1 or 2 along with the regularity conditions 3 and 4. Suppose \mathcal{F} is a collection of balls which cover E in the sense of Definition 9.12.4. Then there exists a sequence of disjoint balls, $\{B_i\} \subseteq \mathcal{F}$ such that $\bar{\mu} \left(E \setminus \bigcup_{j=1}^N B_j \right) = 0$ for $N \leq \infty$.

Proof: Let $\mathbf{x} \in E$. Thus \mathbf{x} is the center of arbitrarily small balls from \mathcal{F} . Since μ is finite on compact sets, only countably many can fail to have $\mu(\partial B(\mathbf{x}, r)) = 0$. Leave the balls out which have $\mu(\partial B(\mathbf{x}, r)) > 0$. Let \mathcal{F}' denote the closures of the balls of \mathcal{F}' . Thus, for these balls, $\mu(\partial B(\mathbf{x}, r)) = 0$. Since for each $\mathbf{x} \in E$ there are only countably many exceptions, \mathcal{F}' is still a Vitali cover of E . Therefore, by Corollary 9.12.3 there is a disjoint sequence of these balls of \mathcal{F}' , $\{\bar{B}_i\}_{i=1}^\infty$ for which $\bar{\mu} \left(E \setminus \bigcup_{j=1}^N \bar{B}_j \right) = 0$. However, since their boundaries have μ measure zero, it follows $\bar{\mu} \left(E \setminus \bigcup_{j=1}^N B_j \right) = 0$, $N \leq \infty$. \blacksquare

9.13 Differentiation of Increasing Functions

As a spectacular application of the covering theorem, is the famous theorem that an increasing function has a derivative a.e. Here the a.e. refers to Lebesgue measure, the Stieltjes measure from the increasing function $F(x) = x$.

Definition 9.13.1 The Dini derivatives are as follows. In these formulas, f is a real

valued function defined on \mathbb{R} .

$$\begin{aligned} D^+f(x) &\equiv \lim_{r \rightarrow 0+} \left(\sup_{0 < u \leq r} \frac{f(x+u) - f(x)}{u} \right) \equiv \lim_{u \rightarrow 0+} \sup \frac{f(x+u) - f(x)}{u}, \\ D_+f(x) &\equiv \lim_{r \rightarrow 0+} \left(\inf_{0 < u \leq r} \frac{f(x+u) - f(x)}{u} \right) \equiv \lim_{u \rightarrow 0+} \inf \frac{f(x+u) - f(x)}{u}, \\ D^-f(x) &\equiv \lim_{r \rightarrow 0+} \left(\sup_{0 < u \leq r} \frac{f(x) - f(x-u)}{u} \right) \equiv \lim_{u \rightarrow 0+} \sup \frac{f(x) - f(x-u)}{u}, \\ D_-f(x) &\equiv \lim_{r \rightarrow 0+} \left(\inf_{0 < u \leq r} \frac{f(x) - f(x-u)}{u} \right) \equiv \lim_{u \rightarrow 0+} \inf \frac{f(x) - f(x-u)}{u}. \end{aligned}$$

Lemma 9.13.2 *The function $f : \mathbb{R} \rightarrow \mathbb{R}$ has a derivative if and only if all the Dini derivatives are equal.*

Proof: If $D^+f(x) = D_+f(x)$, then if u is small enough, let y_n be a decreasing sequence converging to x . Then

$$0 = D^+f(x) - D_+f(x) \geq \limsup_{n \rightarrow \infty} \frac{f(y_n) - f(x)}{y_n - x} - \liminf_{n \rightarrow \infty} \frac{f(y_n) - f(x)}{y_n - x}$$

and so the limit of the difference quotient exists for any such $\{y_n\}$. Thus the derivative from the right exists at x . Therefore, $D^+f(x) > D_+f(x)$ if and only if there is no right derivative. Similarly $D^-f(x) > D_-f(x)$ if and only if there is no derivative from the left at x . Also, there is a derivative if and only if there is a derivative from the left, right and the two are equal. This happens when $D^+f(x) = D_-f(x) = D^-f(x) = D_+f(x)$. Thus this happens if and only if all Dini derivatives are equal. ■

The Lebesgue measure of single points is 0 and so we do not need to worry about whether the intervals are closed in using Corollary 9.12.3.

Let $\Delta f(I) = f(b) - f(x)$ or $f(x) - f(a)$ if I is an interval having end points $a < b$ with x the midpoint. Now suppose $\{J_j\}$ are disjoint intervals contained in I . Then, since f is increasing, $\Delta f(I) \geq \sum_j \Delta f(J_j)$. In this notation, the above lemma implies that if $D^-f(x) > b$ or $D^+f(x) > b$, then for each $\varepsilon > 0$ there is an interval J of length less than ε which is centered at x and $\frac{\Delta f(J)}{(1/2)m(J)} > b$ where $m(J)$ is the Lebesgue measure of J which is the length of J . If either $D_-f(x)$ or $D_+f(x) < a$, the above lemma implies that for each $\varepsilon > 0$ there exists I centered at x with $|I| < \varepsilon$ and $\frac{\Delta f(I)}{(1/2)m(I)} < a$. For example, if $D^-f(x) < a$, there exists a sequence $y_n \uparrow x$ with

$$\frac{f(y_n) - f(x)}{y_n - x} = \frac{f(x) - f(y_n)}{x - y_n} < a$$

so let I_n be the interval centered at x which has left end point y_n .

Note that the set of jumps J of an increasing function must be countable because these jumps determine disjoint open intervals of the form $(f(x-), f(x+))$ for $x \in J$ and each must contain a rational number of which, there are only countably many.

Lemma 9.13.3 *An increasing function f is Borel measurable and its derivatives are Borel measurable functions.*

Proof: The set of jumps J is countable so it is a Borel set of measure zero, an F_σ set. Since f is increasing, its only points of discontinuity are points where it has a jump. Hence it is continuous off this set J . $f^{-1}([c, \infty))$ is an interval of the form $[d, \infty)$ or (d, ∞) . Thus f is Borel measurable. Consider D^+f for $x \notin J$. Let $g_r(x) \equiv \sup_{0 < u \leq r} \frac{f(x+u) - f(x)}{u}$. Thus it is the supremum of functions continuous on J^C . Hence if $x \in g_r^{-1}(c, \infty)$, $c \geq 0$, and $x \notin J^C$, $\frac{f(x+u) - f(x)}{u} > c$ for some $0 < u \leq r$. It follows that, since f is continuous at x , if \hat{x} is close enough to x , it is also true that $\frac{f(\hat{x}+u) - f(\hat{x})}{u} > c$. Thus if $x \in g_r^{-1}(c, \infty) \cap J^C$, then for some δ_x small enough, $(x - \delta_x, x + \delta_x) \subseteq g_r^{-1}(c, \infty)$. Hence, $g_r^{-1}(c, \infty) \cap J^C$ is the intersection of an open set, the union of the intervals $(x - \delta_x, x + \delta_x)$ for $x \in g_r^{-1}(c, \infty) \cap J^C$, with J^C a Borel set. It follows that g_r is decreasing in r and is measurable because, since J is countable, $[g_r > c]$ is the union of a countable set with a Borel set. Thus for all x ,

$$D^+f(x) = \lim_{r \rightarrow 0^+} \left(\sup_{0 < u \leq r} \frac{f(x+u) - f(x)}{u} \right) = \lim_{r \rightarrow 0^+} (g_r(x)) = \lim_{r_n \rightarrow 0} g_{r_n}(x)$$

where r_n is a decreasing sequence converging to 0. It follows that $x \rightarrow D^+f(x)$ is Borel measurable as claimed because it is the limit of Borel measurable functions. Similar reasoning shows that the other derivatives are measurable also. ■

Theorem 9.13.4 *Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be increasing. Then $f'(x)$ exists for all x off a set of measure zero.*

Proof: Let N_{ab} for $0 < a < b$ denote either

$$\{x : D^+f(x) > b > a > D_+f(x)\}, \{x : D^-f(x) > b > a > D_-f(x)\},$$

or

$$\{x : D^-f(x) > b > a > D_+f(x)\}, \{x : D^+f(x) > b > a > D_-f(x)\}$$

From the above lemma, N_{ab} is measurable. Assume that N_{ab} is bounded and let V be open with $V \supseteq N_{ab}$, $m(N_{ab}) + \varepsilon > m(V)$. By Corollary 9.12.3 and the above discussion, there are open, disjoint intervals $\{I_n\}$, each centered at a point of N_{ab} such that

$$\frac{2\Delta f(I_n)}{m(I_n)} < a, \quad m(N_{ab}) = m(N_{ab} \cap \cup_i I_i) = \sum_i m(N_{ab} \cap I_i)$$

Now do for $N_{ab} \cap I_i$ what was just done for N_{ab} and get disjoint intervals J_i^j contained in I_i with

$$\frac{2\Delta f(J_i^j)}{m(J_i^j)} > b, \quad m(N_{ab} \cap I_i) = \sum_j m(N_{ab} \cap I_i \cap J_i^j)$$

Then

$$\begin{aligned} a(m(N_{ab}) + \varepsilon) &> am(V) \geq a \sum_i m(I_i) > \sum_i 2\Delta f(I_i) \geq \sum_i \sum_j 2\Delta f(J_i^j) \\ &\geq b \sum_i \sum_j m(J_i^j) \geq b \sum_i \sum_j m(J_i^j \cap N_{ab}) = b \sum_i m(N_{ab} \cap I_i) = bm(N_{ab}) \end{aligned}$$

Since ε is arbitrary and $a < b$, this shows $m(N_{ab}) = 0$. If N_{ab} is not bounded, apply the above to $N_{ab} \cap (-r, r)$ and conclude this has measure 0. Hence so does N_{ab} .

The countable union of N_{ab} for a, b positive rational and N_{ab} defined in any of the above ways is an exceptional set off which $D^+f(x) = D_+f(x) \geq D^-f(x) \geq D_-f(x) \geq D^+f(x)$ and so these are all equal. This shows that off a set of measure zero, the function has a derivative a.e. ■

9.14 Exercises

1. Suppose you have (X, \mathcal{F}, μ) where $\mathcal{F} \supseteq \mathcal{B}(X)$ and also $\mu(B(x_0, r)) < \infty$ for all $r > 0$. Let $S(x_0, r) \equiv \{x \in X : d(x, x_0) = r\}$. Show that

$$\{r > 0 : \mu(S(x_0, r)) > 0\}$$

cannot be uncountable. Explain why there exists a strictly increasing sequence $r_n \rightarrow \infty$ such that $\mu(x : d(x, x_0) = r_n) = 0$. In other words, the skin of the ball has measure zero except for possibly countably many values of the radius r .

2. Lebesgue measure was discussed. Recall that $m((a, b)) = b - a$ and it is defined on a σ algebra which contains the Borel sets, more generally on $\mathcal{P}(\mathbb{R})$. Also recall that m is translation invariant. Let $x \sim y$ if and only if $x - y \in \mathbb{Q}$. Show this is an equivalence relation. Now let W be a set of positive measure which is contained in $(0, 1)$. For $x \in W$, let $[x]$ denote those $y \in W$ such that $x \sim y$. Thus the equivalence classes partition W . Use axiom of choice to obtain a set $S \subseteq W$ such that S consists of exactly one element from each equivalence class. Let \mathbb{T} denote the rational numbers in $[-1, 1]$. Consider $\mathbb{T} + S \subseteq [-1, 2]$. Explain why $\mathbb{T} + S \supseteq W$. For $\mathbb{T} \equiv \{r_j\}$, explain why the sets $\{r_j + S\}_j$ are disjoint. Now suppose S is measurable. Then show that you have a contradiction if $m(S) = 0$ since $m(W) > 0$ and you also have a contradiction if $m(S) > 0$ because $\mathbb{T} + S$ consists of countably many disjoint sets. Explain why S cannot be measurable. Thus there exists $T \subseteq \mathbb{R}$ such that $m(T) < m(T \cap S) + m(T \cap S^C)$. Is there an open interval (a, b) such that if $T = (a, b)$, then the above inequality holds?
3. Consider the following nested sequence of compact sets, $\{P_n\}$. Let $P_1 = [0, 1]$, $P_2 = [0, \frac{1}{3}] \cup [\frac{2}{3}, 1]$, etc. To go from P_n to P_{n+1} , delete the open interval which is the middle third of each closed interval in P_n . Let $P = \bigcap_{n=1}^{\infty} P_n$. By the finite intersection property of compact sets, $P \neq \emptyset$. Show $m(P) = 0$. If you feel ambitious also show there is a one to one onto mapping of $[0, 1]$ to P . The set P is called the Cantor set. Thus, although P has measure zero, it has the same number of points in it as $[0, 1]$ in the sense that there is a one to one and onto mapping from one to the other. **Hint:** There are various ways of doing this last part but the most enlightenment is obtained by exploiting the topological properties of the Cantor set rather than some silly representation in terms of sums of powers of two and three. All you need to do is use the Schroder Bernstein theorem and show there is an onto map from the Cantor set to $[0, 1]$. If you do this right and remember the theorems about characterizations of compact metric spaces, Proposition 3.5.8 on Page 78, you may get a pretty good idea why every compact metric space is the continuous image of the Cantor set.
4. Consider the sequence of functions defined in the following way. Let $f_1(x) = x$ on $[0, 1]$. To get from f_n to f_{n+1} , let $f_{n+1} = f_n$ on all intervals where f_n is constant. If f_n is nonconstant on $[a, b]$, let $f_{n+1}(a) = f_n(a)$, $f_{n+1}(b) = f_n(b)$, f_{n+1} is piecewise linear and equal to $\frac{1}{2}(f_n(a) + f_n(b))$ on the middle third of $[a, b]$. Sketch a few of

these and you will see the pattern. The process of modifying a nonconstant section of the graph of this function is illustrated in the following picture.



Show $\{f_n\}$ converges uniformly on $[0, 1]$. If $f(x) = \lim_{n \rightarrow \infty} f_n(x)$, show that $f(0) = 0$, $f(1) = 1$, f is continuous, and $f'(x) = 0$ for all $x \notin P$ where P is the Cantor set of Problem 3. This function is called the Cantor function. It is a very important example to remember. Note it has derivative equal to zero a.e. and yet it succeeds in climbing from 0 to 1. Explain why this interesting function cannot be recovered by integrating its derivative. (It is not absolutely continuous, explained later.) **Hint:** This isn't too hard if you focus on getting a careful estimate on the difference between two successive functions in the list considering only a typical small interval in which the change takes place. The above picture should be helpful.

5. \uparrow This problem gives a very interesting example found in the book by McShane [40]. Let $g(x) = x + f(x)$ where f is the strange function of Problem 4. Let P be the Cantor set of Problem 3. Let $[0, 1] \setminus P = \bigcup_{j=1}^{\infty} I_j$ where I_j is open and $I_j \cap I_k = \emptyset$ if $j \neq k$. These intervals are the connected components of the complement of the Cantor set. Show $m(g(I_j)) = m(I_j)$ so $m(g(\bigcup_{j=1}^{\infty} I_j)) = \sum_{j=1}^{\infty} m(g(I_j)) = \sum_{j=1}^{\infty} m(I_j) = 1$. Thus $m(g(P)) = 1$ because $g([0, 1]) = [0, 2]$. By Problem 2 there exists a set, $A \subseteq g(P)$ which is non measurable. Define $\phi(x) = \mathcal{R}_A(g(x))$. Thus $\phi(x) = 0$ unless $x \in P$. Tell why ϕ is measurable. (Recall $m(P) = 0$ and Lebesgue measure is complete.) Now show that $\mathcal{R}_A(y) = \phi(g^{-1}(y))$ for $y \in [0, 2]$. Tell why g is strictly increasing and g^{-1} is continuous but $\phi \circ g^{-1}$ is not measurable. (This is an example of measurable \circ continuous \neq measurable.) Show there exist Lebesgue measurable sets which are not Borel measurable. **Hint:** The function, ϕ is Lebesgue measurable. Now recall that Borel \circ measurable = measurable.
6. Show that every countable set of real numbers is of Lebesgue measure zero.
7. Review the Cantor set in Problem 12 on Page 176. You deleted middle third open intervals. Show that you can take out open intervals in the middle which are not necessarily middle thirds, and end up with a set C which has Lebesgue measure equal to $1 - \varepsilon$. Also show if you can that there exists a continuous and one to one map $f : C \rightarrow J$ where J is the usual Cantor set of Problem 12 which also has measure 0.
8. Recall that every bounded variation function is the difference of two increasing functions. Show that every bounded variation function has a derivative a.e. For a discussion of these, see Definition 11.15.1 on Page 348 below if you have not seen it already.
9. Suppose you have a π system \mathcal{K} of sets of Ω and suppose $\mathcal{G} \supseteq \mathcal{K}$ and that \mathcal{G} is closed with respect to complements and that whenever $\{F_k\}$ is a decreasing sequence of sets of \mathcal{G} it follows that $\bigcap_k F_k \in \mathcal{G}$. Show that then \mathcal{G} contains $\sigma(\mathcal{K})$. This is an alternative formulation of Dynkin's lemma. It was shown after the Dynkin lemma that closure with respect to countable intersections is equivalent.

10. For $x \in \mathbb{R}^p$ to be in $\prod_{i=1}^p A_i$, it means that the i^{th} component of x , x_i is in A_i for each i . Now for $\prod_{i=1}^p (a_i, b_i) \equiv R$, let $V(R) = \prod_{i=1}^p (b_i - a_i)$. Next, for $A \in \mathcal{P}(\mathbb{R}^p)$ let

$$\mu(A) \equiv \inf \left\{ \sum_k V(R^k) : A \subseteq \cup_k R^k \right\}$$

This is just like one dimensional Lebesgue measure except that instead of open intervals, we are using open boxes R^k . Show the following.

- (a) μ is an outer measure.
- (b) $\mu(\prod_{i=1}^p [a_i, b_i]) = \prod_{i=1}^p (b_i - a_i) = \mu(\prod_{i=1}^p (a_i, b_i))$.
- (c) If $\text{dist}(A, B) > 0$, then $\mu(A) + \mu(B) = \mu(A \cup B)$ so $\mathcal{B}(\mathbb{R}^p) \subseteq \mathcal{F}$ the set of sets measurable with respect to this outer measure μ .

This is Lebesgue measure on \mathbb{R}^p . **Hint:** Suppose for some $j, b_j - a_j < \varepsilon$. Show that $\mu(\prod_{i=1}^p (a_i, b_i)) \leq \varepsilon \prod_{i \neq j} (b_i - a_i)$. Now use this to show that if you have a covering by finitely many open boxes, such that the sum of their volumes is less than some number, you can replace with a covering of open boxes which also has the sum of their volumes less than that number but which has each box with sides less than δ . To do this, you might consider replacing each box in the covering with 2^{mp} open boxes obtained by bisecting each side m times where m is small enough that each little box has sides smaller than $\delta/2$ in each of the finitely many boxes in the cover and then fatten each of these just a little to cover up what got left out and retain the sum of the volumes of the little boxes to still be less than the number you had.

11. \uparrow Show that Lebesgue measure defined in the above problem is both inner and outer regular and is translation invariant.
12. Let $(\Omega, \mathcal{F}, \mu)$ be a measure space and let $s(\omega) = \sum_{i=0}^n c_i \mathcal{X}_{E_i}(\omega)$ where the E_i are distinct measurable sets but the c_i might not be. Thus the c_i are the finitely many values of s . Say each $c_i \geq 0$ and $c_0 = 0$. Define $\int s d\mu$ as $\sum_i c_i \mu(E_i)$. Show that this is well defined and that if you have $s(\omega) = \sum_{i=1}^n c_i \mathcal{X}_{E_i}(\omega), t(\omega) = \sum_{j=1}^m d_j \mathcal{X}_{F_j}(\omega)$, then for a, b nonnegative numbers, $as(\omega) + bt(\omega)$ can be written also in this form and that $\int (as + bt) d\mu = a \int s d\mu + b \int t d\mu$. **Hint:** $s(\omega) = \sum_i \sum_j c_i \mathcal{X}_{E_i \cap F_j}(\omega) = \sum_j \sum_i c_i \mathcal{X}_{E_i \cap F_j}(\omega)$ and $(as + bt)(\omega) = \sum_j \sum_i (ac_i + bd_j) \mathcal{X}_{E_i \cap F_j}(\omega)$.
13. \uparrow Having defined the integral of nonnegative simple functions in the above problem, letting f be nonnegative and measurable. Define

$$\int f d\mu \equiv \sup \left\{ \int s d\mu : 0 \leq s \leq f, s \text{ simple} \right\}.$$

Show that if f_n is nonnegative and measurable and $n \rightarrow f_n(\omega)$ is increasing, show that for $f(\omega) = \lim_{n \rightarrow \infty} f_n(\omega)$, it follows that $\int f d\mu = \lim_{n \rightarrow \infty} \int f_n d\mu$. **Hint:** Show $\int f_n d\mu$ is increasing to something $\alpha \leq \infty$. Explain why $\int f d\mu \geq \alpha$. Now pick a nonnegative simple function $s \leq f$. For $r \in (0, 1), [f_n > rs] \equiv E_n$ is increasing in n and $\cup_n E_n = \Omega$. Tell why $\int f_n d\mu \geq \int \mathcal{X}_{E_n} f_n d\mu \geq r \int s d\mu$. Let $n \rightarrow \infty$ and show that $\alpha \geq r \int s d\mu$. Now explain why $\alpha \geq r \int f d\mu$. Since r is arbitrary, $\alpha \geq \int f d\mu \geq \alpha$.

14. \uparrow Show that if f, g are nonnegative and measurable and $a, b \geq 0$, then

$$\int (af + bg) d\mu = a \int f d\mu + b \int g d\mu$$

9.15 Multifunctions and Their Measurability

This is an introduction to the idea of measurable multifunctions. This is a very important topic which has surprising usefulness in nonlinear analysis and other areas and not enough attention is paid to it. As an application, I will give a proof of Kuratowski's theorem and also an interesting fixed point result in which the fixed point is a measurable function of ω in a measure space. One of the main references for this material is the book Papageorgiu and Hu [31] where you can find more of this kind of thing.

9.15.1 The General Case

Let X be a separable complete metric space and let (Ω, \mathcal{F}) be a set and a σ algebra of subsets of Ω . A multifunction, is a map from Ω to the nonempty subsets of X . Thus Γ is a multifunction if for each ω , $\Gamma(\omega) \neq \emptyset$. For more on the theorems presented in this section, see [31].

Definition 9.15.1 Define $\Gamma^-(S) \equiv \{\omega \in \Omega : \Gamma(\omega) \cap S \neq \emptyset\}$. When

$$\Gamma^-(U) \in \mathcal{F}$$

for all U open, we say that Γ is measurable.

More can be said than what follows, but the following is the essential idea for a measurable multifunction.

Theorem 9.15.2 The following are equivalent for any measurable space consisting only of a set Ω and a σ algebra \mathcal{F} . Here nothing is known about $\Gamma(\omega)$ other than that is a nonempty set.

1. For all U open in X , $\Gamma^-(U) \in \mathcal{F}$ where $\Gamma^-(U) \equiv \{\omega : \Gamma(\omega) \cap U \neq \emptyset\}$
2. There exists a sequence, $\{\sigma_n\}$ of measurable functions satisfying $\sigma_n(\omega) \in \Gamma(\omega)$ such that for all $\omega \in \Omega$, $\Gamma(\omega) = \{\sigma_n(\omega) : n \in \mathbb{N}\}$. These functions are called measurable selections.

Proof: First 1.) \Rightarrow 2.). A measurable selection will be obtained in $\overline{\Gamma(\omega)}$. Let $D \equiv \{x_n\}_{n=1}^\infty$ be a countable dense subset of X . For $\omega \in \Omega$, let $\psi_1(\omega) = x_n$ where n is the smallest integer such that $\Gamma(\omega) \cap B(x_n, 1) \neq \emptyset$. Therefore, $\psi_1(\omega)$ has countably many values, x_{n_1}, x_{n_2}, \dots where $n_1 < n_2 < \dots$. Now the set on which ψ_1 has the value x_n is as follows: $\{\omega : \psi_1 = x_n\} =$

$$\{\omega : \Gamma(\omega) \cap B(x_n, 1) \neq \emptyset\} \cap [\Omega \setminus \bigcup_{k < n} \{\omega : \Gamma(\omega) \cap B(x_k, 1) \neq \emptyset\}] \in \mathcal{F}.$$

Thus ψ_1 is measurable and $\text{dist}(\psi_1(\omega), \Gamma(\omega)) < 1$. Let $\Omega_n \equiv \{\omega \in \Omega : \psi_1(\omega) = x_n\}$. Then $\Omega_n \in \mathcal{F}$ and $\Omega_n \cap \Omega_m = \emptyset$ for $n \neq m$ and $\bigcup_{n=1}^\infty \Omega_n = \Omega$ because if ω is given, $\Gamma(\omega)$ does intersect some $B(x_n, 1)$. Let

$$D_n \equiv \{x_k \in D : x_k \in B(x_n, 1)\}.$$

Now for each n , and $\omega \in \Omega_n$, let $\psi_2(\omega) = x_k$ where k is the smallest index such that $x_k \in D_n$ and $B(x_k, \frac{1}{2}) \cap \Gamma(\omega) \neq \emptyset$. Thus

$$\text{dist}(\psi_2(\omega), \Gamma(\omega)) < \frac{1}{2}, \quad d(\psi_2(\omega), \psi_1(\omega)) < 1. \quad (9.25)$$

This defines $\psi_2(\omega)$ on Ω_n and so it defines ψ_2 on Ω satisfying 9.25. Continue this way, obtaining ψ_k a measurable function such that

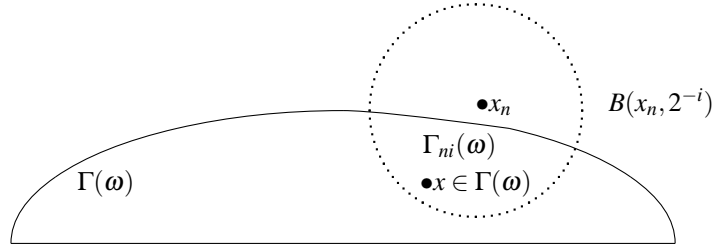
$$\text{dist}(\psi_k(\omega), \Gamma(\omega)) < \frac{1}{2^{k-1}}, \quad d(\psi_{k+1}(\omega), \psi_k(\omega)) < \frac{1}{2^{k-2}}.$$

Then for each ω , $\{\psi_k(\omega)\}$ is a Cauchy sequence of measurable functions converging to a point, $\sigma(\omega) \in \overline{\Gamma(\omega)}$. This has shown that if Γ is measurable, there exists a measurable selection, $\sigma(\omega) \in \overline{\Gamma(\omega)}$. Of course, if $\Gamma(\omega)$ is closed, then $\sigma(\omega) \in \Gamma(\omega)$. Note that this had nothing to do with any measure.

It remains to show that there exists a sequence of these measurable selections σ_n such that the conclusion of 2.) holds. To do this define for a single $\omega \in \Omega$

$$\Gamma_{ni}(\omega) \equiv \begin{cases} \Gamma(\omega) \cap B(x_n, 2^{-i}) & \text{if } \Gamma(\omega) \cap B(x_n, 2^{-i}) \neq \emptyset \\ \Gamma(\omega) & \text{otherwise when there is empty intersection} \end{cases}.$$

The following picture illustrates $\Gamma_{ni}(\omega)$ when ω is such that there is nonempty intersection. Also, given $x \in \Gamma(\omega)$, and i , there is x_n from the countable dense set such that the situation of the picture occurs.



Is Γ_{ni} measurable? If so, then from the above, it has a measurable selection σ_{ni} and the set of these σ_{ni} must have the property that $\{\sigma_{ni}(\omega)\}_{n,i}$ is dense in $\Gamma(\omega)$ for each ω .

Let U be open. Then

$$\begin{aligned} \{\omega : \Gamma_{ni}(\omega) \cap U \neq \emptyset\} &= \{\omega : \Gamma(\omega) \cap B(x_n, 2^{-i}) \cap U \neq \emptyset\} \cup \\ &\quad [\{\omega : \Gamma(\omega) \cap B(x_n, 2^{-i}) = \emptyset\} \cap \{\omega : \Gamma(\omega) \cap U \neq \emptyset\}] \\ &= \{\omega : \Gamma(\omega) \cap B(x_n, 2^{-i}) \cap U \neq \emptyset\} \cup \\ &\quad [(\Omega \setminus \{\omega : \Gamma(\omega) \cap B(x_n, 2^{-i}) \neq \emptyset\}) \cap \{\omega : \Gamma(\omega) \cap U \neq \emptyset\}], \end{aligned}$$

a measurable set. Thus Γ_{ni} is measurable as hoped.

By what was just shown, there exists σ_{ni} , a measurable function such that $\sigma_{ni}(\omega) \in \overline{\Gamma_{ni}(\omega)} \subseteq \overline{\Gamma(\omega)}$ for all $\omega \in \Omega$. If $x \in \overline{\Gamma(\omega)}$, then $x \in \overline{B(x_n, 2^{-(i+1)})}$ whenever x_n is close enough to x . Thus both $x, \sigma_{n(i+1)}(\omega)$ are in $\overline{B(x_n, 2^{-(i+1)})}$ and so $|\sigma_{n(i+1)}(\omega) - x| < 2^{-i}$. It follows that condition 2.) holds with the countable dense subset of $\overline{\Gamma(\omega)}$ being the $\{\sigma_{ni}(\omega)\}$. Note that this had nothing to do with a measure.

Now consider why 2.) \Rightarrow 1.). We have $\{\sigma_n(\omega)\} \subseteq \Gamma(\omega)$ and σ_n is measurable and $\bigcup_n \sigma_n(\omega)$ equals $\overline{\Gamma(\omega)}$. Why is Γ a measurable multifunction? Let U be an open set

$$\begin{aligned} \Gamma^-(U) &\equiv \{\omega : \Gamma(\omega) \cap U \neq \emptyset\} = \{\omega : \overline{\Gamma(\omega)} \cap U \neq \emptyset\} \\ &= \bigcup_n \sigma_n^{-1}(U) \in \mathcal{F} \quad \blacksquare \end{aligned}$$

For much more on multi-functions, you should see the book by Hu and Papageorgiou. [31] The above proof follows the presentation in this book but there is more to be seen there where complete measures are included in the theory and an equivalence is shown between strong measurability, about to be discussed, and measurability without an assumption that the multifunction has compact values.

9.15.2 A Special Case When $\Gamma(\omega)$ Compact

Measurability is a statement that $\Gamma^-(U) \in \mathcal{F}$ whenever U is open.

Definition 9.15.3 *A multifunction Γ is strongly measurable if $\Gamma^-(F) \in \mathcal{F}$ for all F closed.*

Observation 9.15.4 *If Γ is strongly measurable, then it is measurable because if you have U open in a metric space, it is the countable union of closed sets F_n . Hence $\Gamma^-(U) = \cup_k \Gamma^-(F_k) \in \mathcal{F}$.*

Now suppose $\Gamma(\omega)$ is compact for every ω and that $\Gamma^-(U) \in \mathcal{F}$ for every U open. Then let F be a closed set and let $\{U_n\}$ be a decreasing sequence of open sets whose intersection equals F such that also, for all n , $U_n \supseteq \overline{U_{n+1}}$. Then

$$\Gamma(\omega) \cap F = \cap_n \Gamma(\omega) \cap U_n = \cap_n \Gamma(\omega) \cap \overline{U_n}$$

Now because of compactness, the set on the left is nonempty if and only if each set $\Gamma(\omega) \cap \overline{U_n}$ on the right is also nonempty. Thus $\Gamma^-(F) = \cap_n \Gamma^-(U_n) \in \mathcal{F}$. It follows that in this special case, the two conditions, measurability and strong measurability are equivalent. Note that there is no condition on measures or completeness or any such thing. This proves the following proposition.

Proposition 9.15.5 *Let X be a Polish space and let $\Gamma : X \rightarrow \mathcal{P}(X)$ have compact values. Then Γ is measurable if and only if it is strongly measurable, the latter being the statement that $\Gamma^-(C)$ is measurable whenever C is closed.*

Recall how if a function f is measurable, then $f^{-1}(\text{Borel set}) \in \mathcal{F}$. Something like this happens in case Γ is strongly measurable. Let Γ be strongly measurable. Let \mathcal{G} be the sets G such that $\Gamma^-(G)$ and $\Gamma^-(G^C)$ are both in \mathcal{F} . Then clearly \mathcal{G} is closed with respect to complements. If $G \in \mathcal{G}$ is G^C ? Is $\Gamma^-(G^C)$ and $\Gamma^-((G^C)^C)$ in \mathcal{F} ? This is just the definition of what it means to be in \mathcal{G} . Also if you have $\{G_i\} \subseteq \mathcal{G}$, Then

$$\Gamma^-(\cup_i G_i) = \cup_i \Gamma^-(G_i) \in \mathcal{F}$$

and so \mathcal{G} is closed with respect to countable unions. Hence \mathcal{G} must contain the Borel sets because it is a σ algebra and the closed sets are in \mathcal{G} . Thus $\Gamma^-(G) \in \mathcal{F}$ whenever G is Borel.

9.15.3 Kuratowski's Theorem

Also there is a useful corollary from Theorem 9.15.2 and Proposition 9.15.5.

Corollary 9.15.6 *Let $K(\omega)$ be a compact subset of a separable metric space X and suppose $\{u_j(\omega)\}_{j=1}^\infty \subseteq K(\omega)$ with each $\omega \rightarrow u_j(\omega)$ measurable into X . Then there exists $u(\omega) \in K(\omega)$ such that $\omega \rightarrow u(\omega)$ is measurable into X and a subsequence $n(\omega)$ depending on ω such that $\lim_{n(\omega) \rightarrow \infty} u_{n(\omega)}(\omega) = u(\omega)$.*

Proof: Define $\Gamma_n(\omega) = \overline{\cup_{k \geq n} u_k(\omega)}$. This is a nonempty compact subset of $K(\omega) \subseteq X$. I claim that $\omega \rightarrow \Gamma_n(\omega)$ is a measurable multifunction into X . It is necessary to show that $\Gamma_n^-(O)$ defined as $\{\omega : \Gamma_n(\omega) \cap O \neq \emptyset\}$ is measurable whenever O is open in X . For $\omega \in \Gamma_n^-(O)$ it means that some $u_k(\omega) \in O, k \geq n$. Thus $\Gamma_n^-(O) = \cup_{k \geq n} u_k^{-1}(O)$ and this is measurable by the assumption that each u_k is. Since $\Gamma_n^-(\omega)$ is compact, it is also strongly measurable by Proposition 9.15.5, meaning that $\Gamma^-(H)$ is measurable whenever H is closed. Now, let $\Gamma(\omega)$ be defined as $\Gamma(\omega) \equiv \cap_n \Gamma_n(\omega)$ and then for H closed, $\Gamma^-(H)$ is nonempty if and only if $\Gamma_n^-(H)$ is nonempty for each n and $\Gamma^-(H) = \cap_n \Gamma_n^-(H)$ and each set in the intersection is measurable, so this shows that $\omega \rightarrow \Gamma(\omega)$ is also (strongly) measurable. Therefore, it has a measurable selection $u(\omega)$. It follows from the definition of $\Gamma(\omega)$ that there exists a subsequence $n(\omega)$ such that $u(\omega) = \lim_{n(\omega) \rightarrow \infty} u_{n(\omega)}(\omega)$. ■

This corollary makes possible a fairly short proof of the very amazing and enormously significant Kuratowski theorem [35] which gives measurability of maximums of Carathéodory functions.

Definition 9.15.7 *The functions $f : \Omega \times E \rightarrow \mathbb{R}$ in which $f(\cdot, \omega)$ is continuous and $\omega \rightarrow f(x, \omega)$ is measurable are called Carathéodory functions.*

Now here is the Kuratowski theorem.

Theorem 9.15.8 *Let E be a compact metric space and let (Ω, \mathcal{F}) be a measure space. Suppose $\psi : E \times \Omega \rightarrow \mathbb{R}$ has the property that $x \rightarrow \psi(x, \omega)$ is continuous and $\omega \rightarrow \psi(x, \omega)$ is measurable. Then there exists a measurable function f having values in E such that $\psi(f(\omega), \omega) = \max_{x \in E} \psi(x, \omega)$. Furthermore, $\omega \rightarrow \psi(f(\omega), \omega)$ is measurable.*

Proof: Let $C = \{e_i\}_{i=1}^\infty$ be a countable dense subset of E . For example, take the union of $1/2^n$ nets for all n . Let $C_n \equiv \{e_1, \dots, e_n\}$. Let $\omega \rightarrow f_n(\omega)$ be measurable and satisfy $\psi(f_n(\omega), \omega) = \sup_{x \in C_n} \psi(x, \omega)$. This is easily done as follows. Let

$$B_k \equiv \{\omega : \psi(e_k, \omega) \geq \psi(e_j, \omega) \text{ for all } j \neq k\}.$$

Then let $A_1 \equiv B_1$ and if A_1, \dots, A_k have been chosen, let $A_{k+1} \equiv B_{k+1} \setminus \left(\bigcup_{j=1}^k B_j\right)$. Thus each A_k is measurable and you let $f_n(\omega) \equiv e_k$ for $\omega \in A_k$. Using Corollary 9.15.6, there is measurable $f(\omega)$ and a subsequence $n(\omega) \geq n$ such that $f_{n(\omega)}(\omega) \rightarrow f(\omega)$. Then by continuity, $\psi(f(\omega), \omega) = \lim_{n(\omega) \rightarrow \infty} \psi(f_{n(\omega)}(\omega), \omega)$ and this is an increasing sequence in this limit. Hence $\psi(f(\omega), \omega) \geq \sup_{x \in C_n} \psi(x, \omega)$ for each n and so $\psi(f(\omega), \omega) \geq \sup_{x \in C} \psi(x, \omega) = \sup_{x \in E} \psi(x, \omega)$. Since f is measurable, it is the limit of a sequence $\{g_n(\omega)\}$ such that g_n has finitely many values occurring on measurable sets, Theorem 9.1.7. Hence, by continuity, $\psi(f(\omega), \omega) = \lim_{n \rightarrow \infty} \psi(g_n(\omega), \omega)$ and since $\omega \rightarrow \psi(g_n(\omega), \omega)$ is measurable, so is $\psi(f(\omega), \omega)$. ■

One can generalize fairly easily. It is the same argument but carrying around more ω .

Theorem 9.15.9 *Let $E(\omega)$ be a compact metric space in a separable metric space (X, d) and suppose that $\omega \rightarrow E(\omega)$ is a measurable multifunction where (Ω, \mathcal{F}) be a*

measure space. Suppose $\psi_\omega : E(\omega) \times \Omega \rightarrow \mathbb{R}$ has the property that $x \rightarrow \psi_\omega(x, \omega)$ is continuous and $\omega \rightarrow \psi_\omega(x(\omega), \omega)$ is measurable if $x(\omega) \in E(\omega)$ and $\omega \rightarrow x(\omega)$ is measurable ($x(\omega)$ a measurable selection of $E(\omega)$). Then there exists a measurable function f with $f(\omega) \in E(\omega)$ such that $\psi_\omega(f(\omega), \omega) = \max_{x \in E(\omega)} \psi_\omega(x, \omega)$. Furthermore, $\omega \rightarrow \psi_\omega(f(\omega), \omega)$ is measurable.

Proof: Let $C(\omega) = \{e_i(\omega)\}_{i=1}^\infty$ be a countable dense subset of $E(\omega)$ with each $e_i(\omega)$ measurable. This countable dense subset exists by Theorem 9.15.2. Let

$$C_n(\omega) \equiv \{e_1(\omega), \dots, e_n(\omega)\}.$$

Let $\omega \rightarrow f_n(\omega)$ be measurable and satisfy

$$\psi_\omega(f_n(\omega), \omega) = \sup_{x \in C_n(\omega)} \psi_\omega(x, \omega).$$

This is easily done as follows. Let

$$B_k \equiv \{\omega : \psi_\omega(e_k(\omega), \omega) \geq \psi_\omega(e_j(\omega), \omega) \text{ for all } j \neq k\}.$$

Then let $A_1 \equiv B_1$ and if A_1, \dots, A_k have been chosen, let $A_{k+1} \equiv B_{k+1} \setminus \left(\bigcup_{j=1}^k B_j\right)$. Thus each A_k is measurable, and you let $f_n(\omega) \equiv e_k(\omega)$ for $\omega \in A_k$, so $f_n(\omega) \in E(\omega)$ and f_n is measurable. Using Corollary 9.15.6, there is measurable $f(\omega)$ and a subsequence $n(\omega) \geq n$ such that $f_{n(\omega)}(\omega) \rightarrow f(\omega)$. Then by continuity,

$$\psi_\omega(f(\omega), \omega) = \lim_{n(\omega) \rightarrow \infty} \psi_\omega(f_{n(\omega)}(\omega), \omega)$$

and this is an increasing sequence in this limit. Hence

$$\psi_\omega(f(\omega), \omega) \geq \sup_{x \in C_n(\omega)} \psi_\omega(x, \omega)$$

for each n and so

$$\psi_\omega(f(\omega), \omega) \geq \sup_{x \in C(\omega)} \psi_\omega(x, \omega) = \sup_{x \in E(\omega)} \psi_\omega(x, \omega).$$

Since f is measurable, it follows by assumption, that $\omega \rightarrow \psi_\omega(f(\omega), \omega)$ is measurable. ■

Note the following: If you have the simpler situation where $\psi(x, \omega)$ defined on $X \times \Omega$ with $x \rightarrow \psi(x, \omega)$ continuous and $\omega \rightarrow \psi(x, \omega)$ measurable but $E(\omega)$ a compact measurable multifunction as above, then the conditions will hold because you would have $\omega \rightarrow \psi(x(\omega), \omega)$ is measurable if $x(\omega)$ is. Indeed, $x(\omega)$ is the limit of a sequence $\{x_n(\omega)\}$ such that x_n has finitely many values on measurable sets, Theorem 9.1.7. Hence, by continuity, $\psi(x(\omega), \omega) = \lim_{n \rightarrow \infty} \psi(x_n(\omega), \omega)$ and since $\omega \rightarrow \psi(x_n(\omega), \omega)$ is measurable, so is $\psi(x(\omega), \omega)$.

9.15.4 Measurability of Fixed Points

As an interesting application is a consideration of the existence of measurable Brouwer fixed points. This is really quite amazing since Brouwer fixed points are not obtained as the limit of a sequence of iterates although the above Sperner's lemma algorithm provides an

algorithm for finding one. This is a very nice application of the marvelous Kuratowski theorem. It is possible to get this result directly from Corollary 9.15.6 applied to the Sperner's lemma method of proving the Brouwer fixed point theorem. We sent in a paper once which did this and we thought the result was amazing. However, I had forgotten about Kuratowski's theorem which makes this very easy. Fortunately, the referee knew this theorem.

Theorem 9.15.10 *Let K be a closed convex bounded subset of \mathbb{R}^p . Let*

$$x \rightarrow f(x, \omega) : K \rightarrow K$$

be continuous for each ω and $\omega \rightarrow f(x, \omega)$ is measurable, meaning inverse images of sets open in K are in \mathcal{F} where (Ω, \mathcal{F}) is a measurable space. Then there exists $x(\omega) \in K$ such that $\omega \rightarrow x(\omega)$ is measurable and $f(x(\omega), \omega) = x(\omega)$.

Proof: Simply consider $E = K$ and $\psi(x, \omega) \equiv -|x - f(x, \omega)|$. It has a maximum $x(\omega)$ for each ω thanks to continuity of $f(\cdot, \omega)$. Thanks to the Brouwer fixed point theorem, this $x(\omega)$ must be a fixed point. By the above Kuratowski theorem, one of these $x(\omega)$ is measurable. Obviously, by continuity of $f(\cdot, \omega)$, $\omega \rightarrow f(x(\omega), \omega)$ is measurable. ■

If desired, you can extend this to the case where $K(\omega)$ is a measurable multifunction.

9.15.5 Other Measurability Considerations

Here are some other general considerations about measurable multifunctions. The first has to do with getting a new measurable multifunction from old ones and the second has to do with measurability of ε nets. These are technical results which are sometimes useful, for example, if you want to generalize to the Schauder fixed point theorem.

Lemma 9.15.11 *Suppose $f : K(\omega) \times \Omega \rightarrow X, K \subseteq X$. Here X is Polish space, separable complete metric space, and (Ω, \mathcal{F}) is a measurable space. Also $\omega \rightarrow K(\omega)$ is a measurable multifunction as in Theorem 9.15.2. Also suppose*

1. $\omega \rightarrow f(x, \omega)$ is measurable and $x \rightarrow f(x, \omega)$ is continuous.
2. $\mathcal{K}(\omega) \equiv f(K(\omega), \omega)$.

Then you can conclude that $\omega \rightarrow \mathcal{K}(\omega)$ is a measurable multifunction. If $\mathcal{K}(\omega)$ is compact, then it is also strongly measurable.

Proof: Let $\{x_n(\omega)\}$ be a countable dense subset of $K(\omega)$, each x_n measurable. Then if U is open,

$$\{\omega : \mathcal{K}(\omega) \cap U \neq \emptyset\} = \bigcup_{n=1}^{\infty} f(x_n(\cdot), \cdot)^{-1}(U) \quad (9.26)$$

and each of the sets in the union is measurable. The latter claim follows from the continuity of $f(\cdot, \omega)$. If $x(\omega)$ is measurable, then we can express it as the limit of functions s_n which have finitely many values on measurable sets for which $\omega \rightarrow f(s_n(\omega), \omega)$ is clearly measurable. Then $f(x(\omega), \omega)$ is the limit of $f(s_n(\omega), \omega)$. The reason for the equality in 9.26 is as follows. It is clear that the right side is contained in the left. Now if $\mathcal{K}(\omega) \cap U \neq \emptyset$, then by definition, $f(x, \omega) \in U$ for some $x \in K(\omega)$ but then by continuity, $f(x_n(\omega), \omega) \in U$ also for some $x_n(\omega)$ close to x . Thus the two sets are actually equal. Thus $\omega \rightarrow \mathcal{K}(\omega)$ is measurable. If $\mathcal{K}(\omega)$ has compact values it will be strongly measurable as discussed in Proposition 9.15.5. ■

There is also the following general result about the existence of a measurable ε net. This is formulated in Banach space because it is convenient to add. A Banach space is just a complete normed vector space. It could also be formulated in Polish space with a little more difficulty. One just defines things a little differently.

Proposition 9.15.12 *Let $\omega \rightarrow \mathcal{K}(\omega)$ be a measurable multifunction where $\mathcal{K}(\omega)$ is a pre compact set. Recall this means its closure is compact. Thus $\mathcal{K}(\omega)$ must have an ε net for each $\varepsilon > 0$. Then for each $\varepsilon > 0$, there exists $N(\omega)$ and measurable functions $y_j, j = 1, 2, \dots, N(\omega)$, $y_j(\omega) \in \mathcal{K}(\omega)$, such that $\bigcup_{j=1}^{N(\omega)} B(y_j(\omega), \varepsilon) \supseteq \mathcal{K}(\omega)$ for each ω . Also $\omega \rightarrow N(\omega)$ is measurable.*

Proof: Suppose that $\omega \rightarrow \mathcal{K}(\omega)$ is a measurable multifunction having compact values in X a Banach space. Let $\{\sigma_n(\omega)\}$ be the measurable selections such that for each ω , $\{\sigma_n(\omega)\}_{n=1}^\infty$ is dense in $\mathcal{K}(\omega)$. Let $y_1(\omega) \equiv \sigma_1(\omega)$. Now let $2(\omega)$ be the first index larger than 1 such that $\|\sigma_{2(\omega)}(\omega) - \sigma_1(\omega)\| > \frac{\varepsilon}{2}$. Thus $2(\omega) = k$ on the measurable set

$$\left\{ \omega \in \Omega : \|\sigma_k(\omega) - \sigma_1(\omega)\| > \frac{\varepsilon}{2} \right\} \cap \left\{ \omega \in \Omega : \bigcap_{j=1}^{k-1} \|\sigma_j(\omega) - \sigma_1(\omega)\| \leq \frac{\varepsilon}{2} \right\}$$

Suppose $1(\omega), 2(\omega), \dots, (m-1)(\omega)$ have been chosen such that this is a strictly increasing sequence for each ω , each is a measurable function, and for $i, j \leq m-1$,

$$\|\sigma_{i(\omega)}(\omega) - \sigma_{j(\omega)}(\omega)\| > \frac{\varepsilon}{2}.$$

Each $\omega \rightarrow \sigma_{j(\omega)}(\omega)$ is measurable since it equals $\sum_{k=1}^\infty \mathcal{X}_{[i(\omega)=k]}(\omega) \sigma_k(\omega)$. Then $m(\omega)$ will be the first index larger than $(m-1)(\omega)$ such that

$$\|\sigma_{m(\omega)}(\omega) - \sigma_{j(\omega)}(\omega)\| > \frac{\varepsilon}{2}$$

for all $j(\omega) < m(\omega)$. Thus $\omega \rightarrow m(\omega)$ is also measurable because it equals k on the measurable set

$$\left(\bigcap \left\{ \omega : \|\sigma_k(\omega) - \sigma_{j(\omega)}(\omega)\| > \frac{\varepsilon}{2}, j \leq m-1 \right\} \right) \cap \{ \omega : (m-1)(\omega) < k \}$$

$$\cap \left(\bigcup \left\{ \omega : \|\sigma_{k-1}(\omega) - \sigma_{j(\omega)}(\omega)\| \leq \frac{\varepsilon}{2}, j \leq m-1 \right\} \right)$$

The top line says that it does what is wanted and the second says it is the first after $(m-1)(\omega)$ which does so.

Since $\mathcal{K}(\omega)$ is a pre compact set, it follows that the above measurable set will be empty for all $m(\omega)$ sufficiently large called $N(\omega)$, also a measurable function, and so the process ends. Let $y_i(\omega) \equiv \sigma_{i(\omega)}(\omega)$. Then this gives the desired measurable ε net. The fact that

$$\bigcup_{i=1}^{N(\omega)} B(y_i(\omega), \varepsilon) \supseteq \mathcal{K}(\omega)$$

follows because if there exists $z \in \mathcal{K}(\omega) \setminus \left(\bigcup_{i=1}^{N(\omega)} B(y_i(\omega), \varepsilon) \right)$, then $B(z, \frac{\varepsilon}{2})$ would have empty intersection with all of the balls $B(y_i(\omega), \frac{\varepsilon}{3})$ and by density of the $\sigma_i(\omega)$ in $\mathcal{K}(\omega)$, there would be some $\sigma_l(\omega)$ contained in $B(z, \frac{\varepsilon}{3})$ for arbitrarily large l and so the process would not have ended as shown above. ■

9.16 Exercises

1. Using some form of Kuratowski's theorem show the following: Let $K(\omega)$ be a closed convex bounded subset of \mathbb{R}^n where $\omega \rightarrow K(\omega)$ is a measurable multifunction. Let $x \rightarrow f(x, \omega) : K(\omega) \rightarrow K(\omega)$ be continuous for each ω and $\omega \rightarrow f(x, \omega)$ is measurable, meaning inverse images of sets open in \mathbb{R}^n are in \mathcal{F} where (Ω, \mathcal{F}) is a measurable space. Then there exists $x(\omega) \in K(\omega)$ such that $\omega \rightarrow x(\omega)$ is measurable and $f(x(\omega), \omega) = x(\omega)$.
2. If you have $K(\omega)$ a closed convex nonempty set in \mathbb{R}^n and also $\omega \rightarrow K(\omega)$ is a measurable multifunction, show $\omega \rightarrow P_{K(\omega)}x$ is measurable where $P_{K(\omega)}$ is the projection map which gives the closest point in $K(\omega)$. Consider Corollary 6.3.2 on Page 163 or Theorem 11.6.8 and Problem 10 on Page 152 to see the use of this projection map. Also you may want to use Theorem 9.15.2 involving the countable dense subset of $K(\omega)$ consisting of measurable functions.
3. Let $\omega \rightarrow K(\omega)$ be a measurable multifunction in \mathbb{R}^p and let $K(\omega)$ be convex, closed, and compact for each ω . Let $A(\cdot, \omega) : K(\omega) \rightarrow \mathbb{R}^p$ be continuous and $\omega \rightarrow A(x, \omega)$ be measurable. Then if $\omega \rightarrow y(\omega)$ is measurable, there exists measurable $\omega \rightarrow x(\omega)$ such that for all $z \in K(\omega)$,

$$(y(\omega) - A(x(\omega), \omega), z(\omega) - x(\omega)) \leq 0$$

This is a measurable version of Browder's lemma, a very important result in nonlinear analysis. **Hint:** You want to have for each ω ,

$$P_{K(\omega)}(y(\omega) - A(x, \omega) + x) = x$$

Use Problem 2 and the measurability of Brouwer fixed points discussed above.

4. In the situation of the above problem, suppose also that $\lim_{|x| \rightarrow \infty} \frac{(A(x, \omega), x)}{|x|} = \infty$. Show that there exists measurable $x(\omega)$ such that $A(x(\omega), \omega) = y(\omega)$. **Hint:** Let $x_n(\omega)$ be the solution of Problem 3 in which $K_n = B(0, n)$. Show that these are bounded for each ω . Then use Corollary 9.15.6 to get $x(\omega)$, a suitable limit such that $A(x(\omega), \omega) = y(\omega)$.

Chapter 10

The Abstract Lebesgue Integral

The general Lebesgue integral requires a measure space, $(\Omega, \mathcal{F}, \mu)$ and, to begin with, a nonnegative measurable function. I will use Lemma 2.5.3 about interchanging two supremums frequently. Also, I will use the observation that if $\{a_n\}$ is an increasing sequence of points of $[0, \infty]$, then $\sup_n a_n = \lim_{n \rightarrow \infty} a_n$ which is obvious from the definition of sup.

10.1 Nonnegative Measurable Functions

10.1.1 Riemann Integrals for Decreasing Functions

First of all, the notation $[g < f]$ means $\{\omega \in \Omega : g(\omega) < f(\omega)\}$ with other variants of this notation being similar. Also, the convention, $0 \cdot \infty = 0$ will be used to simplify the presentation whenever it is convenient to do so. The notation $a \wedge b$ means the minimum of a and b .

Definition 10.1.1 Let $f : [a, b] \rightarrow [0, \infty]$ be decreasing. Note that ∞ is a possible value. Define

$$\int_a^b f(\lambda) d\lambda \equiv \lim_{M \rightarrow \infty} \int_a^b M \wedge f(\lambda) d\lambda = \sup_M \int_a^b M \wedge f(\lambda) d\lambda$$

where $a \wedge b$ means the minimum of a and b . Note that for f bounded,

$$\sup_M \int_a^b M \wedge f(\lambda) d\lambda = \int_a^b f(\lambda) d\lambda$$

where the integral on the right is the usual Riemann integral because eventually $M > f$. For f a nonnegative decreasing function defined on $[0, \infty)$,

$$\int_0^\infty f d\lambda \equiv \lim_{R \rightarrow \infty} \int_0^R f d\lambda = \sup_{R > 1} \int_0^R f d\lambda = \sup_R \sup_{M > 0} \int_0^R f \wedge M d\lambda$$

Since decreasing bounded functions are Riemann integrable, the above definition is well defined. For a discussion of this, see Calculus of One and Many Variables on the web site or any elementary Calculus text. Now here is an obvious property.

Lemma 10.1.2 Let f be a decreasing nonnegative function defined on an interval $[a, b]$. Then if $[a, b] = \cup_{k=1}^m I_k$ where $I_k \equiv [a_k, b_k]$ and the intervals I_k are non overlapping, it follows

$$\int_a^b f d\lambda = \sum_{k=1}^m \int_{a_k}^{b_k} f d\lambda.$$

Proof: This follows from the computation,

$$\int_a^b f d\lambda \equiv \lim_{M \rightarrow \infty} \int_a^b f \wedge M d\lambda = \lim_{M \rightarrow \infty} \sum_{k=1}^m \int_{a_k}^{b_k} f \wedge M d\lambda = \sum_{k=1}^m \int_{a_k}^{b_k} f d\lambda$$

Note both sides could equal $+\infty$. ■

In all considerations below, we assume h is fairly small, certainly much smaller than R . Thus $R - h > 0$.

Lemma 10.1.3 *Let g be a decreasing nonnegative function defined on an interval $[0, R]$. Then*

$$\int_0^R g \wedge M d\lambda = \sup_{h>0} \sum_{i=1}^{m(R,h)} (g(ih) \wedge M) h$$

where $m(h, R) \in \mathbb{N}$ satisfies $R - h < hm(h, R) \leq R$.

Proof: Since $g \wedge M$ is a decreasing bounded function the lower sums converge to the integral as $h \rightarrow 0$. Thus

$$\int_0^R g \wedge M d\lambda = \lim_{h \rightarrow 0} \left(\sum_{i=1}^{m(R,h)} (g(ih) \wedge M) h + (g(R) \wedge M) (R - hm(h, R)) \right)$$

Now the last term in the above is no more than Mh and so the above is

$$\lim_{h \rightarrow 0} \left(\sum_{i=1}^{m(R,h)} (g(ih) \wedge M) h \right) = \sup_{h>0} \left(\sum_{i=1}^{m(R,h)} (g(ih) \wedge M) h \right). \blacksquare$$

10.1.2 The Lebesgue Integral for Nonnegative Functions

Here is the definition of the Lebesgue integral of a function which is measurable and has values in $[0, \infty]$.

Definition 10.1.4 *Let $(\Omega, \mathcal{F}, \mu)$ be a measure space and suppose $f : \Omega \rightarrow [0, \infty]$ is measurable. Then define $\int f d\mu \equiv \int_0^\infty \mu([f > \lambda]) d\lambda$ which makes sense because $\lambda \rightarrow \mu([f > \lambda])$ is nonnegative and decreasing.*

Note that if $f \leq g$, then $\int f d\mu \leq \int g d\mu$ because $\mu([f > \lambda]) \leq \mu([g > \lambda])$.
For convenience $\sum_{i=1}^0 a_i \equiv 0$.

Lemma 10.1.5 *In the above definition, $\int f d\mu = \sup_{h>0} \sum_{i=1}^\infty \mu([f > hi]) h$*

Proof: Let $m(h, R) \in \mathbb{N}$ satisfy $R - h < hm(h, R) \leq R$. Then $\lim_{R \rightarrow \infty} m(h, R) = \infty$ and so from Lemma 10.1.3,

$$\begin{aligned} \int f d\mu &\equiv \int_0^\infty \mu([f > \lambda]) d\lambda = \sup_M \sup_R \int_0^R \mu([f > \lambda]) \wedge M d\lambda \\ &= \sup_M \sup_{R>0} \sup_{h>0} \sum_{k=1}^{m(h,R)} (\mu([f > kh]) \wedge M) h \end{aligned}$$

Hence, switching the order of the sups, this equals

$$\begin{aligned} \sup_{R>0} \sup_{h>0} \sup_M \sum_{k=1}^{m(h,R)} (\mu([f > kh]) \wedge M) h &= \sup_{R>0} \sup_{h>0} \lim_{M \rightarrow \infty} \sum_{k=1}^{m(h,R)} (\mu([f > kh]) \wedge M) h \\ &= \sup_{h>0} \sup_R \sum_{k=1}^{m(R,h)} (\mu([f > kh])) h = \sup_{h>0} \sum_{k=1}^\infty (\mu([f > kh])) h. \blacksquare \end{aligned}$$

10.2 Nonnegative Simple Functions

To begin with, here is a useful lemma.

Lemma 10.2.1 *If $f(\lambda) = 0$ for all $\lambda > a$, where f is a decreasing nonnegative function, then $\int_0^\infty f(\lambda) d\lambda = \int_0^a f(\lambda) d\lambda$.*

Proof: From the definition,

$$\begin{aligned} \int_0^\infty f(\lambda) d\lambda &= \lim_{R \rightarrow \infty} \int_0^R f(\lambda) d\lambda = \sup_{R > 1} \int_0^R f(\lambda) d\lambda = \sup_{R > 1} \sup_M \int_0^R f(\lambda) \wedge M d\lambda \\ &= \sup_M \sup_{R > 1} \int_0^R f(\lambda) \wedge M d\lambda = \sup_M \sup_{R > 1} \int_0^a f(\lambda) \wedge M d\lambda \\ &= \sup_M \int_0^a f(\lambda) \wedge M d\lambda \equiv \int_0^a f(\lambda) d\lambda. \blacksquare \end{aligned}$$

Now the Lebesgue integral for a nonnegative function has been defined, what does it do to a nonnegative simple function? Recall a nonnegative simple function is one which has finitely many nonnegative real values which it assumes on measurable sets. Thus a simple function can be written in the form $s(\omega) = \sum_{i=1}^n c_i \mathcal{X}_{E_i}(\omega)$ where the c_i are each nonnegative, the distinct values of s .

Lemma 10.2.2 *Let $s(\omega) = \sum_{i=1}^p a_i \mathcal{X}_{E_i}(\omega)$ be a nonnegative simple function where the E_i are distinct but the a_i might not be. Thus the values of s are the a_i . Then*

$$\int s d\mu = \sum_{i=1}^p a_i \mu(E_i). \quad (10.1)$$

Proof: Without loss of generality, assume $0 \equiv a_0 < a_1 \leq a_2 \leq \dots \leq a_p$ and that $\mu(E_i) < \infty, i > 0$. Here is why. If $\mu(E_i) = \infty$, then letting $a \in (a_{i-1}, a_i)$, by Lemma 10.2.1, the left side is

$$\begin{aligned} \int_0^{a_p} \mu([s > \lambda]) d\lambda &\geq \int_{a_0}^{a_i} \mu([s > \lambda]) d\lambda \\ &\equiv \sup_M \int_0^{a_i} \mu([s > \lambda]) \wedge M d\lambda \geq \sup_M \sup_M M \mu(E_i) a_i = \infty \end{aligned}$$

and so both sides of 10.1 are equal to ∞ . Thus it can be assumed for each $i, \mu(E_i) < \infty$. Then it follows from Lemma 10.2.1 and Lemma 10.1.2,

$$\begin{aligned} \int_0^\infty \mu([s > \lambda]) d\lambda &= \int_0^{a_p} \mu([s > \lambda]) d\lambda = \sum_{k=1}^p \int_{a_{k-1}}^{a_k} \mu([s > \lambda]) d\lambda \\ &= \sum_{k=1}^p (a_k - a_{k-1}) \sum_{i=k}^p \mu(E_i) = \sum_{i=1}^p \mu(E_i) \sum_{k=1}^i (a_k - a_{k-1}) = \sum_{i=1}^p a_i \mu(E_i) \blacksquare \end{aligned}$$

Note that this is the same result as in Problem 12 on Page 270 but here there is no question about the definition of the integral of a simple function being well defined.

Lemma 10.2.3 *If $a, b \geq 0$ and if s and t are nonnegative simple functions, then*

$$\int as + bt d\mu = a \int s d\mu + b \int t d\mu.$$

Proof: Let $s(\omega) = \sum_{i=1}^n \alpha_i \chi_{A_i}(\omega)$, $t(\omega) = \sum_{j=1}^m \beta_j \chi_{B_j}(\omega)$ where α_i are the distinct values of s and the β_j are the distinct values of t . Clearly $as + bt$ is a nonnegative simple function because it has finitely many values on measurable sets. In fact, $(as + bt)(\omega) = \sum_{j=1}^m \sum_{i=1}^n (a\alpha_i + b\beta_j) \chi_{A_i \cap B_j}(\omega)$ where the sets $A_i \cap B_j$ are disjoint and measurable. By Lemma 10.2.2,

$$\begin{aligned} & \int as + btd\mu \\ &= \sum_{j=1}^m \sum_{i=1}^n (a\alpha_i + b\beta_j) \mu(A_i \cap B_j) = \sum_{i=1}^n a \sum_{j=1}^m \alpha_i \mu(A_i \cap B_j) + b \sum_{j=1}^m \sum_{i=1}^n \beta_j \mu(A_i \cap B_j) \\ &= a \sum_{i=1}^n \alpha_i \mu(A_i) + b \sum_{j=1}^m \beta_j \mu(B_j) = a \int sd\mu + b \int td\mu. \blacksquare \end{aligned}$$

10.3 The Monotone Convergence Theorem

The following is called the monotone convergence theorem. This theorem and related convergence theorems are the reason for using the Lebesgue integral. If $\lim_{n \rightarrow \infty} f_n(\omega) = f(\omega)$ and f_n is increasing in n , then clearly f is also measurable because

$$f^{-1}((a, \infty]) = \cup_{k=1}^{\infty} f_k^{-1}((a, \infty]) \in \mathcal{F}$$

For a different approach to this, see Problem 12 on Page 270.

Theorem 10.3.1 (Monotone Convergence theorem) Suppose that the function f has all values in $[0, \infty]$ and suppose $\{f_n\}$ is a sequence of nonnegative measurable functions having values in $[0, \infty]$ and satisfying

$$\begin{aligned} \lim_{n \rightarrow \infty} f_n(\omega) &= f(\omega) \text{ for each } \omega. \\ \cdots f_n(\omega) &\leq f_{n+1}(\omega) \cdots \end{aligned}$$

Then f is measurable and $\int f d\mu = \lim_{n \rightarrow \infty} \int f_n d\mu$.

Proof: By Lemma 10.1.5 $\lim_{n \rightarrow \infty} \int f_n d\mu = \sup_n \int f_n d\mu$

$$\begin{aligned} &= \sup_n \sup_{h>0} \sum_{k=1}^{\infty} \mu([f_n > kh]) h = \sup_{h>0} \sup_N \sup_n \sum_{k=1}^N \mu([f_n > kh]) h \\ &= \sup_{h>0} \sup_N \sum_{k=1}^N \mu([f > kh]) h = \sup_{h>0} \sum_{k=1}^{\infty} \mu([f > kh]) h = \int f d\mu. \blacksquare \end{aligned}$$

Note how it was important to have $\int_0^{\infty} [f > \lambda] d\lambda$ in the definition of the integral and **not** $[f \geq \lambda]$. You need to have $[f_n > kh] \uparrow [f > kh]$ so $\mu([f_n > kh]) \rightarrow \mu([f > kh])$. To illustrate what goes wrong without the Lebesgue integral, consider the following example.

Example 10.3.2 Let $\{r_n\}$ denote the rational numbers in $[0, 1]$ and let

$$f_n(t) \equiv \begin{cases} 1 & \text{if } t \notin \{r_1, \dots, r_n\} \\ 0 & \text{otherwise} \end{cases}$$

Then $f_n(t) \uparrow f(t)$ where f is the function which is one on the rationals and zero on the irrationals. Each f_n is Riemann integrable (why?) but f is not Riemann integrable because it is everywhere discontinuous. Also, there is a gap between all upper sums and lower sums. Therefore, you can't write $\int f dx = \lim_{n \rightarrow \infty} \int f_n dx$.

An observation which is typically true related to this type of example is this. If you can choose your functions, you don't need the Lebesgue integral. The Riemann Darboux integral is just fine. It is when you can't choose your functions and they come to you as pointwise limits that you really need the superior Lebesgue integral or at least something more general than the Riemann integral. The Riemann integral is entirely adequate for evaluating the seemingly endless lists of boring problems found in calculus books. It is shown later that the two integrals coincide when the Lebesgue integral is taken with respect to Lebesgue measure and the function being integrated is continuous.

10.4 Other Definitions

To review and summarize the above, if $f \geq 0$ is measurable,

$$\int f d\mu \equiv \int_0^\infty \mu([f > \lambda]) d\lambda \quad (10.2)$$

another way to get the same thing for $\int f d\mu$ is to take an increasing sequence of non-negative simple functions, $\{s_n\}$ with $s_n(\omega) \rightarrow f(\omega)$ and then by monotone convergence theorem, $\int f d\mu = \lim_{n \rightarrow \infty} \int s_n$ where if $s_n(\omega) = \sum_{j=1}^m c_j \mathcal{X}_{E_j}(\omega)$, $\int s_n d\mu = \sum_{j=1}^m c_j \mu(E_j)$. Similarly this also shows that for such nonnegative measurable function,

$$\int f d\mu = \sup \left\{ \int s : 0 \leq s \leq f, s \text{ simple} \right\}.$$

Here is an equivalent definition of the integral of a nonnegative measurable function. The fact it is well defined has been discussed above.

Definition 10.4.1 For s a nonnegative simple function,

$$s(\omega) = \sum_{k=1}^n c_k \mathcal{X}_{E_k}(\omega), \int s = \sum_{k=1}^n c_k \mu(E_k).$$

For f a nonnegative measurable function,

$$\int f d\mu = \sup \left\{ \int s : 0 \leq s \leq f, s \text{ simple} \right\}.$$

10.5 Fatou's Lemma

The next theorem, known as Fatou's lemma is another important theorem which justifies the use of the Lebesgue integral.

Theorem 10.5.1 (Fatou's lemma) Let f_n be a nonnegative measurable function. Let $g(\omega) = \liminf_{n \rightarrow \infty} f_n(\omega)$. Then g is measurable and $\int g d\mu \leq \liminf_{n \rightarrow \infty} \int f_n d\mu$. In other words, $\int (\liminf_{n \rightarrow \infty} f_n) d\mu \leq \liminf_{n \rightarrow \infty} \int f_n d\mu$.

Proof: Let $g_n(\omega) = \inf\{f_k(\omega) : k \geq n\}$. Then

$$g_n^{-1}([a, \infty]) = \cap_{k=n}^\infty f_k^{-1}([a, \infty]) = \left(\cup_{k=n}^\infty f_k^{-1}([a, \infty])^c \right)^c \in \mathcal{F}.$$

Thus g_n is measurable by Lemma 9.1.4. Also $g(\omega) = \lim_{n \rightarrow \infty} g_n(\omega)$ so g is measurable because it is the pointwise limit of measurable functions. Now the functions g_n form an

increasing sequence of nonnegative measurable functions so the monotone convergence theorem applies. This yields

$$\int g d\mu = \lim_{n \rightarrow \infty} \int g_n d\mu \leq \liminf_{n \rightarrow \infty} \int f_n d\mu.$$

The last inequality holding because $\int g_n d\mu \leq \int f_n d\mu$. (Note that it is not known whether $\lim_{n \rightarrow \infty} \int f_n d\mu$ exists.) ■

10.6 The Integral's Righteous Algebraic Desires

The monotone convergence theorem shows the integral wants to be linear. This is the essential content of the next theorem.

Theorem 10.6.1 *Let f, g be nonnegative measurable functions and let a, b be non-negative numbers. Then $af + bg$ is measurable and*

$$\int (af + bg) d\mu = a \int f d\mu + b \int g d\mu. \quad (10.3)$$

Proof: By Theorem 9.1.6 on Page 239 there exist increasing sequences of nonnegative simple functions, $s_n \rightarrow f$ and $t_n \rightarrow g$. Then $af + bg$, being the pointwise limit of the simple functions $as_n + bt_n$, is measurable. Now by the monotone convergence theorem and Lemma 10.2.3,

$$\begin{aligned} \int (af + bg) d\mu &= \lim_{n \rightarrow \infty} \int as_n + bt_n d\mu = \lim_{n \rightarrow \infty} \left(a \int s_n d\mu + b \int t_n d\mu \right) \\ &= a \int f d\mu + b \int g d\mu. \quad \blacksquare \end{aligned}$$

As long as you are allowing functions to take the value $+\infty$, you cannot consider something like $f + (-g)$ and so you can't very well expect a satisfactory statement about the integral being linear until you restrict yourself to functions which have values in a vector space. To be linear, a function must be defined on a vector space. This is discussed next.

10.7 The Lebesgue Integral, L^1

The functions considered here have values in \mathbb{C} , which is a vector space. A function f with values in \mathbb{C} is of the form $f = \operatorname{Re} f + i \operatorname{Im} f$ where $\operatorname{Re} f$ and $\operatorname{Im} f$ are real valued functions. In fact $\operatorname{Re} f = \frac{f + \bar{f}}{2}$, $\operatorname{Im} f = \frac{f - \bar{f}}{2i}$.

Definition 10.7.1 *Let $(\Omega, \mathcal{S}, \mu)$ be a measure space and suppose $f : \Omega \rightarrow \mathbb{C}$. Then f is said to be measurable if both $\operatorname{Re} f$ and $\operatorname{Im} f$ are measurable real valued functions.*

Of course there is another definition of measurability which says that inverse images of open sets are measurable. This is equivalent to this new definition.

Lemma 10.7.2 *Let $f : \Omega \rightarrow \mathbb{C}$. Then f is measurable if and only if $\operatorname{Re} f, \operatorname{Im} f$ are both real valued measurable functions. Also if f, g are complex measurable functions and a, b are complex scalars, then $af + bg$ is also measurable.*

Proof: \Rightarrow Suppose first that f is measurable. Recall that \mathbb{C} is considered as \mathbb{R}^2 with (x, y) being identified with $x + iy$. Thus the open sets of \mathbb{C} can be obtained with either of the two equivalent norms $|z| \equiv \sqrt{(\operatorname{Re} z)^2 + (\operatorname{Im} z)^2}$ or $\|z\|_\infty = \max(\operatorname{Re} z, \operatorname{Im} z)$. Therefore, if f is measurable

$$\operatorname{Re} f^{-1}(a, b) \cap \operatorname{Im} f^{-1}(c, d) = f^{-1}((a, b) + i(c, d)) \in \mathcal{F}$$

In particular, you could let $(c, d) = \mathbb{R}$ and conclude that $\operatorname{Re} f$ is measurable because in this case, the above reduces to the statement that $\operatorname{Re} f^{-1}(a, b) \in \mathcal{F}$. Similarly $\operatorname{Im} f$ is measurable.

\Leftarrow Next, if each of $\operatorname{Re} f$ and $\operatorname{Im} f$ are measurable, then

$$f^{-1}((a, b) + i(c, d)) = \operatorname{Re} f^{-1}(a, b) \cap \operatorname{Im} f^{-1}(c, d) \in \mathcal{F}$$

and so, since every open set is the countable union of sets of the form $(a, b) + i(c, d)$, it follows that f is measurable.

Now consider the last claim. Let $h : \mathbb{C} \times \mathbb{C} \rightarrow \mathbb{C}$ be given by $h(z, w) \equiv az + bw$. Then h is continuous. If f, g are complex valued measurable functions, consider the complex valued function, $h \circ (f, g) : \Omega \rightarrow \mathbb{C}$. Then

$$(h \circ (f, g))^{-1}(\text{open}) = (f, g)^{-1}(h^{-1}(\text{open})) = (f, g)^{-1}(\text{open})$$

Now letting U, V be open in \mathbb{C} , $(f, g)^{-1}(U \times V) = f^{-1}(U) \cap g^{-1}(V) \in \mathcal{F}$. Since every open set in $\mathbb{C} \times \mathbb{C}$ is the countable union of sets of the form $U \times V$, it follows that $(f, g)^{-1}(\text{open})$ is in \mathcal{F} . Thus $af + bg$ is also complex measurable. ■

As is always the case for complex numbers, $|z|^2 = (\operatorname{Re} z)^2 + (\operatorname{Im} z)^2$. Also, for g a real valued function, one can consider its positive and negative parts defined respectively as

$$g^+(x) \equiv \frac{g(x) + |g(x)|}{2}, \quad g^-(x) = \frac{|g(x)| - g(x)}{2}.$$

Thus $|g| = g^+ + g^-$ and $g = g^+ - g^-$ and both g^+ and g^- are measurable nonnegative functions if g is measurable.

Then the following is the definition of what it means for a complex valued function f to be in $L^1(\Omega)$.

Definition 10.7.3 Let $(\Omega, \mathcal{F}, \mu)$ be a measure space. Then a complex valued measurable function f is in $L^1(\Omega)$ if $\int |f| d\mu < \infty$. For a function in $L^1(\Omega)$, the integral is defined as follows.

$$\int f d\mu \equiv \int (\operatorname{Re} f)^+ d\mu - \int (\operatorname{Re} f)^- d\mu + i \left[\int (\operatorname{Im} f)^+ d\mu - \int (\operatorname{Im} f)^- d\mu \right]$$

I will show that with this definition, the integral is linear and well defined. First note that it is clearly well defined because all the above integrals are of nonnegative functions and are each equal to a nonnegative real number because for h equal to any of the functions, $|h| \leq |f|$ and $\int |f| d\mu < \infty$.

Here is a lemma which will make it possible to show the integral is linear.

Lemma 10.7.4 Let g, h, g', h' be nonnegative measurable functions in $L^1(\Omega)$ and suppose that $g - h = g' - h'$. Then $\int g d\mu - \int h d\mu = \int g' d\mu - \int h' d\mu$.

Proof: By assumption, $g + h' = g' + h$. Then from the Lebesgue integral's righteous algebraic desires, Theorem 10.6.1, $\int g d\mu + \int h' d\mu = \int g' d\mu + \int h d\mu$ which implies the claimed result. ■

Lemma 10.7.5 *Let $\text{Re}(L^1(\Omega))$ denote the vector space of real valued functions in $L^1(\Omega)$ where the field of scalars is the real numbers. Then $\int d\mu$ is linear on $\text{Re}(L^1(\Omega))$, the scalars being real numbers.*

Proof: First observe that from the definition of the positive and negative parts of a function, $(f + g)^+ - (f + g)^- = f^+ + g^+ - (f^- + g^-)$ because both sides equal $f + g$. Therefore from Lemma 10.7.4 and the definition, it follows from Theorem 10.6.1 that

$$\begin{aligned} \int f + g d\mu &\equiv \int (f + g)^+ - (f + g)^- d\mu = \int f^+ + g^+ d\mu - \int f^- + g^- d\mu \\ &= \int f^+ d\mu + \int g^+ d\mu - \left(\int f^- d\mu + \int g^- d\mu \right) = \int f d\mu + \int g d\mu. \end{aligned}$$

what about taking out scalars? First note that if a is real and nonnegative, then $(af)^+ = af^+$ and $(af)^- = af^-$ while if $a < 0$, then $(af)^+ = -af^-$ and $(af)^- = -af^+$. These claims follow immediately from the above definitions of positive and negative parts of a function. Thus if $a < 0$ and $f \in L^1(\Omega)$, it follows from Theorem 10.6.1 that

$$\begin{aligned} \int af d\mu &\equiv \int (af)^+ d\mu - \int (af)^- d\mu = \int (-a)f^- d\mu - \int (-a)f^+ d\mu \\ &= -a \int f^- d\mu + a \int f^+ d\mu = a \left(\int f^+ d\mu - \int f^- d\mu \right) \equiv a \int f d\mu. \end{aligned}$$

The case where $a \geq 0$ works out similarly but easier. ■

Now here is the main result.

Theorem 10.7.6 *$\int d\mu$ is linear on $L^1(\Omega)$ and $L^1(\Omega)$ is a complex vector space. If $f \in L^1(\Omega)$, then $\text{Re } f$, $\text{Im } f$, and $|f|$ are all in $L^1(\Omega)$. Furthermore, for $f \in L^1(\Omega)$,*

$$\begin{aligned} \int f d\mu &\equiv \int (\text{Re } f)^+ d\mu - \int (\text{Re } f)^- d\mu + i \left[\int (\text{Im } f)^+ d\mu - \int (\text{Im } f)^- d\mu \right] \\ &\equiv \int \text{Re } f d\mu + i \int \text{Im } f d\mu \end{aligned}$$

and the triangle inequality holds,

$$\left| \int f d\mu \right| \leq \int |f| d\mu. \quad (10.4)$$

Also, for every $f \in L^1(\Omega)$ it follows that for every $\varepsilon > 0$ there exists a simple function s such that $|s| \leq |f|$ and $\int |f - s| d\mu < \varepsilon$.

Proof: First consider the claim that the integral is linear. It was shown above that the integral is linear on $\text{Re}(L^1(\Omega))$. Then letting $a + ib, c + id$ be scalars and f, g functions in $L^1(\Omega)$,

$$(a + ib)f + (c + id)g = (a + ib)(\text{Re } f + i\text{Im } f) + (c + id)(\text{Re } g + i\text{Im } g)$$

$$= c \operatorname{Re}(g) - b \operatorname{Im}(f) - d \operatorname{Im}(g) + a \operatorname{Re}(f) + i(b \operatorname{Re}(f) + c \operatorname{Im}(g) + a \operatorname{Im}(f) + d \operatorname{Re}(g))$$

It follows from the definition that

$$\begin{aligned} \int (a + ib)f + (c + id)gd\mu &= \int (c \operatorname{Re}(g) - b \operatorname{Im}(f) - d \operatorname{Im}(g) + a \operatorname{Re}(f))d\mu \\ &\quad + i \int (b \operatorname{Re}(f) + c \operatorname{Im}(g) + a \operatorname{Im}(f) + d \operatorname{Re}(g))d\mu \end{aligned} \quad (10.5)$$

Also, from the definition,

$$\begin{aligned} (a + ib) \int f d\mu + (c + id) \int g d\mu &= (a + ib) \left(\int \operatorname{Re} f d\mu + i \int \operatorname{Im} f d\mu \right) \\ &\quad + (c + id) \left(\int \operatorname{Re} g d\mu + i \int \operatorname{Im} g d\mu \right) \end{aligned}$$

which equals

$$\begin{aligned} &= a \int \operatorname{Re} f d\mu - b \int \operatorname{Im} f d\mu + ib \int \operatorname{Re} f d\mu + ia \int \operatorname{Im} f d\mu \\ &\quad + c \int \operatorname{Re} g d\mu - d \int \operatorname{Im} g d\mu + id \int \operatorname{Re} g d\mu - d \int \operatorname{Im} g d\mu. \end{aligned}$$

Using Lemma 10.7.5 and collecting terms, it follows that this reduces to 10.5. Thus the integral is linear as claimed.

Consider the claim about approximation with a simple function. Letting h equal any of

$$(\operatorname{Re} f)^+, (\operatorname{Re} f)^-, (\operatorname{Im} f)^+, (\operatorname{Im} f)^-, \quad (10.6)$$

It follows from the monotone convergence theorem and Theorem 9.1.6 on Page 239 there exists a nonnegative simple function $s \leq h$ such that $\int |h - s| d\mu < \frac{\varepsilon}{4}$. Therefore, letting s_1, s_2, s_3, s_4 be such simple functions, approximating respectively the functions listed in 10.6, and $s \equiv s_1 - s_2 + i(s_3 - s_4)$,

$$\begin{aligned} \int |f - s| d\mu &\leq \int |(\operatorname{Re} f)^+ - s_1| d\mu + \int |(\operatorname{Re} f)^- - s_2| d\mu \\ &\quad + \int |(\operatorname{Im} f)^+ - s_3| d\mu + \int |(\operatorname{Im} f)^- - s_4| d\mu < \varepsilon \end{aligned}$$

It is clear from the construction that $|s| \leq |f|$.

What about 10.4? Let $\theta \in \mathbb{C}$ be such that $|\theta| = 1$ and $\theta \int f d\mu = |\int f d\mu|$. Then from what was shown above about the integral being linear,

$$\left| \int f d\mu \right| = \theta \int f d\mu = \int \theta f d\mu = \int \operatorname{Re}(\theta f) d\mu \leq \int |f| d\mu.$$

If $f, g \in L^1(\Omega)$, then it is known that for a, b scalars, it follows that $af + bg$ is measurable. See Lemma 10.7.2. Also $\int |af + bg| d\mu \leq \int |a| |f| + |b| |g| d\mu < \infty$. ■

The following corollary follows from this. The conditions of this corollary are sometimes taken as a definition of what it means for a function f to be in $L^1(\Omega)$.

Corollary 10.7.7 $f \in L^1(\Omega)$ if and only if there exists a sequence of complex simple functions, $\{s_n\}$ such that

$$\begin{aligned} s_n(\omega) &\rightarrow f(\omega) \text{ for all } \omega \in \Omega \\ \lim_{m,n \rightarrow \infty} \int (|s_n - s_m|) d\mu &= 0 \end{aligned} \quad (10.7)$$

When $f \in L^1(\Omega)$,

$$\int f d\mu \equiv \lim_{n \rightarrow \infty} \int s_n. \quad (10.8)$$

Proof: From the above theorem, if $f \in L^1$ there exists a sequence of simple functions $\{s_n\}$ such that

$$\int |f - s_n| d\mu < 1/n, \quad s_n(\omega) \rightarrow f(\omega) \text{ for all } \omega$$

Then $\int |s_n - s_m| d\mu \leq \int |s_n - f| d\mu + \int |f - s_m| d\mu \leq \frac{1}{n} + \frac{1}{m}$.

Next suppose the existence of the approximating sequence of simple functions. Then f is measurable because its real and imaginary parts are the limit of measurable functions. By Fatou's lemma, $\int |f| d\mu \leq \liminf_{n \rightarrow \infty} \int |s_n| d\mu < \infty$ because $|\int |s_n| d\mu - \int |s_m| d\mu| \leq \int |s_n - s_m| d\mu$ which is given to converge to 0. Thus $\{\int |s_n| d\mu\}$ is a Cauchy sequence and is therefore, bounded.

In case $f \in L^1(\Omega)$, letting $\{s_n\}$ be the approximating sequence, Fatou's lemma implies

$$\left| \int f d\mu - \int s_n d\mu \right| \leq \int |f - s_n| d\mu \leq \liminf_{m \rightarrow \infty} \int |s_m - s_n| d\mu < \varepsilon$$

provided n is large enough. Hence 10.8 follows. ■

This is a good time to observe the following fundamental observation which follows from a repeat of the above arguments.

Theorem 10.7.8 Suppose $\Lambda(f) \in [0, \infty]$ for all nonnegative measurable functions and suppose that for $a, b \geq 0$ and f, g nonnegative measurable functions,

$$\Lambda(af + bg) = a\Lambda(f) + b\Lambda(g).$$

In other words, Λ wants to be linear. Then Λ has a unique linear extension to the set of measurable functions $\{f \text{ measurable} : \Lambda(|f|) < \infty\}$, this set being a vector space.

10.8 The Dominated Convergence Theorem

One of the major theorems in this theory is the dominated convergence theorem. Before presenting it, here is a technical lemma about limsup and liminf which is really pretty obvious from the definition.

Lemma 10.8.1 Let $\{a_n\}$ be a sequence in $[-\infty, \infty]$. Then $\lim_{n \rightarrow \infty} a_n$ exists if and only if $\liminf_{n \rightarrow \infty} a_n = \limsup_{n \rightarrow \infty} a_n$ and in this case, the limit equals the common value of these two numbers.

Proof: Suppose first $\lim_{n \rightarrow \infty} a_n = a \in \mathbb{R}$. Letting $\varepsilon > 0$ be given, $a_n \in (a - \varepsilon, a + \varepsilon)$ for all n large enough, say $n \geq N$. Therefore, both $\inf\{a_k : k \geq n\}$ and $\sup\{a_k : k \geq n\}$ are contained in $[a - \varepsilon, a + \varepsilon]$ whenever $n \geq N$. It follows $\limsup_{n \rightarrow \infty} a_n$ and $\liminf_{n \rightarrow \infty} a_n$ are

both in $[a - \varepsilon, a + \varepsilon]$, showing $|\liminf_{n \rightarrow \infty} a_n - \limsup_{n \rightarrow \infty} a_n| < 2\varepsilon$. Since ε is arbitrary, the two must be equal and they both must equal a . Next suppose $\lim_{n \rightarrow \infty} a_n = \infty$. Then if $l \in \mathbb{R}$, there exists N such that for $n \geq N$, $l \leq a_n$ and therefore, for such n , $l \leq \inf\{a_k : k \geq n\} \leq \sup\{a_k : k \geq n\}$ and this shows, since l is arbitrary that $\liminf_{n \rightarrow \infty} a_n = \limsup_{n \rightarrow \infty} a_n = \infty$. The case for $-\infty$ is similar.

Conversely, suppose $\liminf_{n \rightarrow \infty} a_n = \limsup_{n \rightarrow \infty} a_n = a$. Suppose first that $a \in \mathbb{R}$. Then, letting $\varepsilon > 0$ be given, there exists N such that if $n \geq N$, $\sup\{a_k : k \geq n\} - \inf\{a_k : k \geq n\} < \varepsilon$. Therefore, if $k, m > N$, and $a_k > a_m$,

$$|a_k - a_m| = a_k - a_m \leq \sup\{a_k : k \geq n\} - \inf\{a_k : k \geq n\} < \varepsilon$$

showing that $\{a_n\}$ is a Cauchy sequence. Therefore, it converges to $a \in \mathbb{R}$, and as in the first part, the \liminf and \limsup both equal a . If $\liminf_{n \rightarrow \infty} a_n = \limsup_{n \rightarrow \infty} a_n = \infty$, then given $l \in \mathbb{R}$, there exists N such that for $n \geq N$, $\inf_{k \geq n} a_k > l$. Therefore, $\lim_{n \rightarrow \infty} a_n = \infty$. The case for $-\infty$ is similar. ■

Here is the dominated convergence theorem.

Theorem 10.8.2 (*Dominated Convergence theorem*) Let $f_n \in L^1(\Omega)$ and suppose

$$f(\omega) = \lim_{n \rightarrow \infty} f_n(\omega),$$

and there exists a measurable function g , with values in $[0, \infty]$,¹ such that

$$|f_n(\omega)| \leq g(\omega) \text{ and } \int g(\omega) d\mu < \infty.$$

Then $f \in L^1(\Omega)$ and

$$0 = \lim_{n \rightarrow \infty} \int |f_n - f| d\mu = \lim_{n \rightarrow \infty} \left| \int f d\mu - \int f_n d\mu \right|$$

Proof: f is measurable by Theorem 9.1.2. Since $|f| \leq g$, it follows that

$$f \in L^1(\Omega) \text{ and } |f - f_n| \leq 2g.$$

By Fatou's lemma (Theorem 10.5.1),

$$\int 2g d\mu \leq \liminf_{n \rightarrow \infty} \int 2g - |f - f_n| d\mu = \int 2g d\mu - \limsup_{n \rightarrow \infty} \int |f - f_n| d\mu.$$

Subtracting $\int 2g d\mu$, $0 \leq -\limsup_{n \rightarrow \infty} \int |f - f_n| d\mu$. Hence

$$\begin{aligned} 0 &\geq \limsup_{n \rightarrow \infty} \left(\int |f - f_n| d\mu \right) \\ &\geq \liminf_{n \rightarrow \infty} \left(\int |f - f_n| d\mu \right) \geq \liminf_{n \rightarrow \infty} \left| \int f d\mu - \int f_n d\mu \right| \geq 0. \end{aligned}$$

This proves the theorem by Lemma 10.8.1 because the \limsup and \liminf are equal. ■

¹Note that, since g is allowed to have the value ∞ , it is not known that $g \in L^1(\Omega)$.

Corollary 10.8.3 Suppose $f_n \in L^1(\Omega)$ and $f(\omega) = \lim_{n \rightarrow \infty} f_n(\omega)$. Suppose also there exist measurable functions, g_n, g with values in $[0, \infty]$ such that

$$\lim_{n \rightarrow \infty} \int g_n d\mu = \int g d\mu$$

$g_n(\omega) \rightarrow g(\omega)$ μ a.e. and both $\int g_n d\mu$ and $\int g d\mu$ are finite. Also suppose $|f_n(\omega)| \leq g_n(\omega)$. Then $\lim_{n \rightarrow \infty} \int |f - f_n| d\mu = 0$.

Proof: It is just like the above. This time $g + g_n - |f - f_n| \geq 0$ and so by Fatou's lemma,

$$\begin{aligned} \int 2g d\mu - \limsup_{n \rightarrow \infty} \int |f - f_n| d\mu &= \lim_{n \rightarrow \infty} \int (g_n + g) d\mu - \limsup_{n \rightarrow \infty} \int |f - f_n| d\mu \\ &= \liminf_{n \rightarrow \infty} \int (g_n + g) d\mu - \limsup_{n \rightarrow \infty} \int |f - f_n| d\mu \\ &= \liminf_{n \rightarrow \infty} \int ((g_n + g) - |f - f_n|) d\mu \geq \int 2g d\mu \end{aligned}$$

and so $-\limsup_{n \rightarrow \infty} \int |f - f_n| d\mu \geq 0$. Thus

$$\begin{aligned} 0 &\geq \limsup_{n \rightarrow \infty} \left(\int |f - f_n| d\mu \right) \\ &\geq \liminf_{n \rightarrow \infty} \left(\int |f - f_n| d\mu \right) \geq \left| \int f d\mu - \int f_n d\mu \right| \geq 0. \blacksquare \end{aligned}$$

Definition 10.8.4 Let E be a measurable subset of Ω . $\int_E f d\mu \equiv \int f \chi_E d\mu$.

If $L^1(E)$ is written, the σ algebra is defined as $\{E \cap A : A \in \mathcal{F}\}$ and the measure is μ restricted to this smaller σ algebra. Clearly, if $f \in L^1(\Omega)$, then $f \chi_E \in L^1(E)$ and if $f \in L^1(E)$, then letting \tilde{f} be the 0 extension of f off of E , it follows $\tilde{f} \in L^1(\Omega)$.

Another very important observation applies to the case where Ω is also a metric space. In this lemma, $\text{spt}(f)$ denotes the closure of the set on which f is nonzero.

Definition 10.8.5 Let K be a set and let V be an open set containing K . Then the notation $K \prec f \prec V$ means that $f(x) = 1$ for all $x \in K$ and $\text{spt}(f)$ is a compact subset of V . $\text{spt}(f)$ is defined as the closure of the set where f is not zero. It is called the “support” of f . A function $f \in C_c(\Omega)$ for Ω a metric space if f is continuous on Ω and $\text{spt}(f)$ is compact. This $C_c(\Omega)$ is called the continuous functions with compact support.

Recall Lemma 3.12.4. Listed next for convenience.

Lemma 10.8.6 Let Ω be a metric space in which the closed balls are compact and let K be a compact subset of V , an open set. Then there exists a continuous function $f : \Omega \rightarrow [0, 1]$ such that $K \prec f \prec V$.

Theorem 10.8.7 Let $(\Omega, \mathcal{S}, \mu)$ be a regular measure space, meaning that μ is inner and outer regular and $\mu(K) < \infty$ for each compact set K . Suppose also that the conclusion of Lemma 3.12.4 holds. Then for each $\varepsilon > 0$ and $f \in L^1(\Omega)$, there is $g \in C_c(\Omega)$ such that $\int_\Omega |f - g| d\mu < \varepsilon$.

Proof: First consider a measurable set E where $\mu(E) < \infty$. Let $K \subseteq E \subseteq V$ where $\mu(V \setminus K) < \varepsilon$. Now let $K \prec h \prec V$. Then

$$\int |h - \mathcal{X}_E| d\mu \leq \int \mathcal{X}_{V \setminus K} d\mu = \mu(V \setminus K) < \varepsilon. \quad (10.9)$$

By Corollary 10.7.7, there is a sequence of simple functions converging pointwise to f such that for $m, n > N$, $\frac{\varepsilon}{2} > \int (|s_n - s_m|) d\mu$. Then let $n \rightarrow \infty$ and apply the dominated convergence theorem to get $\int |f - s_m| d\mu \leq \frac{\varepsilon}{2}$. However, from 10.9, there is g in $C_c(\Omega)$ such that $\int |s_m - g| d\mu < \frac{\varepsilon}{2}$ and so $\int |f - g| d\mu \leq \int |f - s_m| d\mu + \int |s_m - g| d\mu < \varepsilon$. ■

10.9 Some Important General Theory

10.9.1 Eggoroff's Theorem

You might show that a sequence of measurable real or complex valued functions converges on a measurable set. This is Proposition 9.1.8 above. Eggoroff's theorem says that if the set of points where a sequence of measurable functions converges is all but a set of measure zero, then the sequence almost converges uniformly in a certain sense.

Theorem 10.9.1 (Egoroff) *Let $(\Omega, \mathcal{F}, \mu)$ be a finite measure space, $\mu(\Omega) < \infty$ and let f_n, f be complex valued functions such that $\operatorname{Re} f_n, \operatorname{Im} f_n$ are all measurable and*

$$\lim_{n \rightarrow \infty} f_n(\omega) = f(\omega)$$

for all $\omega \notin E$ where $\mu(E) = 0$. Then for every $\varepsilon > 0$, there exists a set,

$$F \supseteq E, \mu(F) < \varepsilon,$$

such that f_n converges uniformly to f on F^C .

Proof: First suppose $E = \emptyset$ so that convergence is pointwise everywhere. It follows then that $\operatorname{Re} f$ and $\operatorname{Im} f$ are pointwise limits of measurable functions and are therefore measurable. Let $E_{km} = \{\omega \in \Omega : |f_n(\omega) - f(\omega)| \geq 1/m \text{ for some } n > k\}$. Note that

$$|f_n(\omega) - f(\omega)| = \sqrt{(\operatorname{Re} f_n(\omega) - \operatorname{Re} f(\omega))^2 + (\operatorname{Im} f_n(\omega) - \operatorname{Im} f(\omega))^2}$$

and so, $[|f_n - f| \geq \frac{1}{m}]$ is measurable. Hence E_{km} is measurable because

$$E_{km} = \bigcup_{n=k+1}^{\infty} \left[|f_n - f| \geq \frac{1}{m} \right].$$

For fixed m , $\bigcap_{k=1}^{\infty} E_{km} = \emptyset$ because f_n converges to f . Therefore, if $\omega \in \Omega$ there exists k such that if $n > k$, $|f_n(\omega) - f(\omega)| < \frac{1}{m}$ which means $\omega \notin E_{km}$. Note also that $E_{km} \supseteq E_{(k+1)m}$. Since $\mu(E_{1m}) < \infty$, Theorem 9.2.4 on Page 242 implies

$$0 = \mu\left(\bigcap_{k=1}^{\infty} E_{km}\right) = \lim_{k \rightarrow \infty} \mu(E_{km}).$$

Let $k(m)$ be chosen such that $\mu(E_{k(m)m}) < \varepsilon 2^{-m}$ and let $F = \bigcup_{m=1}^{\infty} E_{k(m)m}$. Then $\mu(F) < \varepsilon$ because $\mu(F) \leq \sum_{m=1}^{\infty} \mu(E_{k(m)m}) < \sum_{m=1}^{\infty} \varepsilon 2^{-m} = \varepsilon$.

Now let $\eta > 0$ be given and pick m_0 such that $m_0^{-1} < \eta$. If $\omega \in F^C$, then $\omega \in \bigcap_{m=1}^{\infty} E_{k(m)m}^C$. Hence $\omega \in E_{k(m_0)m_0}^C$ so $|f_n(\omega) - f(\omega)| < 1/m_0 < \eta$ for all $n > k(m_0)$. This holds for all $\omega \in F^C$ and so f_n converges uniformly to f on F^C .

Now if $E \neq \emptyset$, consider $\{\mathcal{X}_{E^C} f_n\}_{n=1}^{\infty}$. Each $\mathcal{X}_{E^C} f_n$ has real and imaginary parts measurable and the sequence converges pointwise to $\mathcal{X}_E f$ everywhere. Therefore, from the first part, there exists a set of measure less than ε , F such that on F^C , $\{\mathcal{X}_{E^C} f_n\}$ converges uniformly to $\mathcal{X}_{E^C} f$. Therefore, on $(E \cup F)^C$, $\{f_n\}$ converges uniformly to f . This proves the theorem. ■

10.9.2 The Vitali Convergence Theorem

The Vitali convergence theorem is a convergence theorem which in the case of a finite measure space is superior to the dominated convergence theorem.

Definition 10.9.2 Let $(\Omega, \mathcal{F}, \mu)$ be a measure space and let $\mathfrak{S} \subseteq L^1(\Omega)$. \mathfrak{S} is uniformly integrable if for every $\varepsilon > 0$ there exists $\delta > 0$ such that for all $f \in \mathfrak{S}$

$$\left| \int_E f d\mu \right| < \varepsilon \text{ whenever } \mu(E) < \delta.$$

Lemma 10.9.3 If \mathfrak{S} is uniformly integrable, then $|\mathfrak{S}| \equiv \{|f| : f \in \mathfrak{S}\}$ is uniformly integrable. Also \mathfrak{S} is uniformly integrable if \mathfrak{S} is finite.

Proof: Let $\varepsilon > 0$ be given and suppose \mathfrak{S} is uniformly integrable. First suppose the functions are real valued. Let δ be such that if $\mu(E) < \delta$, then $|\int_E f d\mu| < \frac{\varepsilon}{2}$ for all $f \in \mathfrak{S}$. Let $\mu(E) < \delta$. Then if $f \in \mathfrak{S}$,

$$\begin{aligned} \int_E |f| d\mu &\leq \int_{E \cap [f \leq 0]} (-f) d\mu + \int_{E \cap [f > 0]} f d\mu = \left| \int_{E \cap [f \leq 0]} f d\mu \right| + \left| \int_{E \cap [f > 0]} f d\mu \right| \\ &< \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon. \end{aligned}$$

In general, if \mathfrak{S} is a uniformly integrable set of complex valued functions, the inequalities,

$$\left| \int_E \operatorname{Re} f d\mu \right| \leq \left| \int_E f d\mu \right|, \quad \left| \int_E \operatorname{Im} f d\mu \right| \leq \left| \int_E f d\mu \right|,$$

imply $\operatorname{Re} \mathfrak{S} \equiv \{\operatorname{Re} f : f \in \mathfrak{S}\}$ and $\operatorname{Im} \mathfrak{S} \equiv \{\operatorname{Im} f : f \in \mathfrak{S}\}$ are also uniformly integrable. Therefore, applying the above result for real valued functions to these sets of functions, it follows $|\mathfrak{S}|$ is uniformly integrable also.

For the last part, it suffices to verify a single function in $L^1(\Omega)$ is uniformly integrable. To do so, note that from the dominated convergence theorem, $\lim_{R \rightarrow \infty} \int_{[|f| > R]} |f| d\mu = 0$. Let $\varepsilon > 0$ be given and choose R large enough that $\int_{[|f| > R]} |f| d\mu < \frac{\varepsilon}{2}$. Now let $\mu(E) < \frac{\varepsilon}{2R}$. Then

$$\begin{aligned} \int_E |f| d\mu &= \int_{E \cap [|f| \leq R]} |f| d\mu + \int_{E \cap [|f| > R]} |f| d\mu \\ &< R\mu(E) + \frac{\varepsilon}{2} < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon. \end{aligned}$$

This proves the lemma. ■

The following gives a nice way to identify a uniformly integrable set of functions.

Lemma 10.9.4 Let \mathfrak{S} be a subset of $L^1(\Omega, \mu)$ where $\mu(\Omega) < \infty$. Let $t \rightarrow h(t)$ be a continuous function which satisfies $\lim_{t \rightarrow \infty} \frac{h(t)}{t} = \infty$. Then \mathfrak{S} is uniformly integrable and bounded in $L^1(\Omega)$ if $\sup \{ \int_{\Omega} h(|f|) d\mu : f \in \mathfrak{S} \} = N < \infty$.

Proof: First I show \mathfrak{S} is bounded in $L^1(\Omega; \mu)$ which means there exists a constant M such that for all $f \in \mathfrak{S}$, $\int_{\Omega} |f| d\mu \leq M$. From the properties of h , there exists R_n such that if $t \geq R_n$, then $h(t) \geq nt$. Therefore, $\int_{\Omega} |f| d\mu = \int_{[|f| \geq R_n]} |f| d\mu + \int_{[|f| < R_n]} |f| d\mu$. Letting $n = 1$, and $f \in \mathfrak{S}$,

$$\begin{aligned} \int_{\Omega} |f| d\mu &= \int_{[|f| \geq R_1]} |f| d\mu + \int_{[|f| < R_1]} |f| d\mu \\ &\leq \int_{[|f| \geq R_1]} h(|f|) d\mu + R_1 \mu([|f| < R_1]) \leq N + R_1 \mu(\Omega) \equiv M. \end{aligned} \quad (10.10)$$

Next let E be a measurable set. Then for every $f \in \mathfrak{S}$, it follows from 10.10

$$\begin{aligned} \int_E |f| d\mu &= \int_{[|f| \geq R_n] \cap E} |f| d\mu + \int_{[|f| < R_n] \cap E} |f| d\mu \\ &\leq \frac{1}{n} \int_{\Omega} |f| d\mu + R_n \mu(E) \leq \frac{M}{n} + R_n \mu(E) \end{aligned} \quad (10.11)$$

Let n be large enough that $M/n < \varepsilon/2$ and then let $\mu(E) < \varepsilon/2R_n$. Then 10.11 is less than $\varepsilon/2 + R_n(\varepsilon/2R_n) = \varepsilon$ ■

Letting $h(t) = t^2$, it follows that if all the functions in \mathfrak{S} are bounded, then the collection of functions is uniformly integrable. Another way to discuss uniform integrability is the following. This other way involving equi-integrability is used a lot in probability.

Definition 10.9.5 Let $(\Omega, \mathcal{F}, \mu)$ be a measure space with $\mu(\Omega) < \infty$. A set $\mathfrak{S} \subseteq L^1(\Omega)$ is said to be equi-integrable if for every $\varepsilon > 0$ there exists $\lambda > 0$ sufficiently large, such that $\int_{[|f| > \lambda]} |f| d\mu < \varepsilon$ for all $f \in \mathfrak{S}$.

Then the relation between this and uniform integrability is as follows.

Proposition 10.9.6 In the context of the above definition, \mathfrak{S} is equi-integrable if and only if it is a bounded subset of $L^1(\Omega)$ which is also uniformly integrable.

Proof: \Rightarrow I need to show \mathfrak{S} is bounded and uniformly integrable. First consider bounded. Choose λ to work for $\varepsilon = 1$. Then for all $f \in \mathfrak{S}$,

$$\int |f| d\mu = \int_{[|f| > \lambda]} |f| d\mu + \int_{[|f| \leq \lambda]} |f| d\mu \leq 1 + \lambda \mu(\Omega)$$

Thus it is bounded. Now let E be a measurable subset of Ω . Let λ go with $\varepsilon/2$ in the definition of equi-integrable. Then for all $f \in \mathfrak{S}$,

$$\int_E |f| d\mu \leq \int_{[|f| > \lambda]} |f| d\mu + \int_{E \cap [|f| \leq \lambda]} |f| d\mu \leq \frac{\varepsilon}{2} + \lambda \mu(E)$$

Then let $\mu(E)$ be small enough that $\lambda \mu(E) < \varepsilon/2$ and this shows uniform integrability.

\Leftarrow I need to verify equi-integrable from bounded and uniformly integrable. Let δ be such that if $\mu(E) < \delta$, then $\int_E |f| d\mu < \varepsilon$ for all $f \in \mathfrak{S}$. If not, then there exists $f_n \in \mathfrak{S}$ with $[|f_n| > n] > \delta$. Thus $\int |f_n| d\mu \geq \int_{[|f_n| > n]} |f_n| d\mu \geq n \mu([|f_n| > n]) > n\delta$ and so \mathfrak{S} is not bounded after all. ■

The following theorem is Vitali's convergence theorem.

Theorem 10.9.7 *Let $\{f_n\}$ be a uniformly integrable set of complex valued functions, $\mu(\Omega) < \infty$, and $f_n(x) \rightarrow f(x)$ a.e. where f is a measurable complex valued function. Then $f \in L^1(\Omega)$ and*

$$\lim_{n \rightarrow \infty} \int_{\Omega} |f_n - f| d\mu = 0. \quad (10.12)$$

Proof: First it will be shown that $f \in L^1(\Omega)$. By uniform integrability, there exists $\delta > 0$ such that if $\mu(E) < \delta$, then $\int_E |f_n| d\mu < 1$ for all n . By Egoroff's theorem, there exists a set E of measure less than δ such that on E^C , $\{f_n\}$ converges uniformly. Therefore, for p large enough, and $n > p$, $\int_{E^C} |f_p - f_n| d\mu < 1$ which implies $\int_{E^C} |f_n| d\mu < 1 + \int_{\Omega} |f_p| d\mu$. Then since there are only finitely many functions, f_n with $n \leq p$, there exists a constant, M_1 such that for all n , $\int_{E^C} |f_n| d\mu < M_1$. But also,

$$\int_{\Omega} |f_m| d\mu = \int_{E^C} |f_m| d\mu + \int_E |f_m| d\mu \leq M_1 + 1 \equiv M.$$

Therefore, by Fatou's lemma, $\int_{\Omega} |f| d\mu \leq \liminf_{n \rightarrow \infty} \int |f_n| d\mu \leq M$, showing that $f \in L^1$ as hoped.

Now $\mathfrak{S} \cup \{f\}$ is uniformly integrable so there exists $\delta_1 > 0$ such that if $\mu(E) < \delta_1$, then $\int_E |g| d\mu < \varepsilon/3$ for all $g \in \mathfrak{S} \cup \{f\}$.

By Egoroff's theorem, there exists a set, F with $\mu(F) < \delta_1$ such that f_n converges uniformly to f on F^C . Therefore, there exists m such that if $n > m$, then $\int_{F^C} |f - f_n| d\mu < \frac{\varepsilon}{3}$. It follows that for $n > m$,

$$\int_{\Omega} |f - f_n| d\mu \leq \int_{F^C} |f - f_n| d\mu + \int_F |f| d\mu + \int_F |f_n| d\mu < \frac{\varepsilon}{3} + \frac{\varepsilon}{3} + \frac{\varepsilon}{3} = \varepsilon,$$

which verifies 10.12. ■

10.10 One Dimensional Lebesgue Stieltjes Integral

Let F be an increasing function defined on \mathbb{R} . Let μ be the Lebesgue Stieltjes measure defined in Theorems 9.9.1 and 9.7.4. The conclusions of these theorems are reviewed here.

Theorem 10.10.1 *Let F be an increasing function defined on \mathbb{R} , an integrator function. There exists a function $\mu : \mathcal{P}(\mathbb{R}) \rightarrow [0, \infty]$ which satisfies the conditions of Theorem 9.7.4 in terms of measures of intervals and the inner and outer regularity properties.*

The Lebesgue integral taken with respect to this measure, is called the Lebesgue Stieltjes integral. Note that any real valued continuous function is measurable with respect to \mathcal{S} . This is because if f is continuous, inverse images of open sets are open and open sets are in \mathcal{S} . Thus f is measurable because $f^{-1}((a, b)) \in \mathcal{S}$. Similarly if f has complex values this argument applied to its real and imaginary parts yields the conclusion that f is measurable. This will be denoted here by $\int f d\mu$ but it is often the case that it is denoted as $\int f dF$.

In the case of most interest, where $F(x) = x$, how does the Lebesgue integral compare with the Riemann integral? The short answer is that if f is Riemann integrable, then it is also Lebesgue interable and the two integrals coincide. It is customary to denote the Lebesgue integral in this context as $\int_a^b f dm$.

Theorem 10.10.2 *Suppose f is Riemann integrable on an interval $[a, b]$. Then f is also Lebesgue integrable and the two integrals are the same.*

Proof: It suffices to consider the case that f is nonnegative. Otherwise, one simply considers the positive and negative parts of the real and imaginary parts of the function. Thus f is a bounded function and there is a decreasing sequence of upper step functions, denoted as $\{u_n\}$ and an increasing sequence of lower step functions denoted as $\{l_n\}$ such that

$$\int_a^b l_n dt \leq \int_a^b f dt \leq \int_a^b u_n dt, \left| \int_a^b u_n dt - \int_a^b l_n dt \right| < 2^{-n}$$

Since f must be bounded, it can be assumed that $|u_n(t)|, |l_n(t)| < M$ for some constant M . Let $g(t) = \lim_{n \rightarrow \infty} u_n(t)$ and $h(t) = \lim_{n \rightarrow \infty} l_n(t)$. Then from the dominated convergence theorem (Why?) one obtains

$$\begin{aligned} \int_a^b f dt &= \lim_{n \rightarrow \infty} \int_a^b l_n dt = \lim_{n \rightarrow \infty} \int_a^b l_n dm = \int_a^b h dm \\ &\leq \int_a^b g dm \leq \lim_{n \rightarrow \infty} \int_a^b u_n dm = \lim_{n \rightarrow \infty} \int_a^b u_n dt = \int_a^b f dt \end{aligned}$$

Also, from the construction, $h(t) \leq f(t) \leq g(t)$. From the above, $\int_a^b |g(t) - h(t)| dm = 0$. It follows that g is measurable (why?) and $f(t) = g(t)$ for m a.e. t . (why?) By completeness of the measure, it follows that f is Lebesgue measurable and $\int_a^b f dm = \int_a^b g dm = \int_a^b h dm = \int_a^b f dt$. (why?) ■

If you have seen the Darboux Stieltjes integral, defined like the Riemann integral in terms of upper and lower sums, the following compares the Lebesgue Stieltjes integral with this one also. For f a continuous function, how does the Lebesgue Stieltjes integral compare with the Darboux Stieltjes integral? To answer this question, here is a technical lemma.

Lemma 10.10.3 *Let D be a countable subset of \mathbb{R} and suppose $a, b \notin D$. Also suppose f is a continuous function defined on $[a, b]$. Then there exists a sequence of functions $\{s_n\}$ of the form $s_n(x) \equiv \sum_{k=1}^{m_n} f(z_{k-1}^n) \mathcal{R}_{[z_{k-1}^n, z_k^n)}(x)$ such that each $z_k^n \notin D$ and*

$$\sup \{|s_n(x) - f(x)| : x \in [a, b]\} < 1/n.$$

Proof: First note that D contains no intervals. To see this let $D = \{d_k\}_{k=1}^\infty$. If D has an interval of length 2ε , let I_k be an interval centered at d_k which has length $\varepsilon/2^k$. Therefore, the sum of the lengths of these intervals is no more than $\sum_{k=1}^\infty \frac{\varepsilon}{2^k} = \varepsilon$. Thus D cannot contain an interval of length 2ε . Since ε is arbitrary, D cannot contain any interval.

Since f is continuous, it follows from Theorem 3.7.4 on Page 82 that f is uniformly continuous. Therefore, there exists $\delta > 0$ such that if $|x - y| \leq 3\delta$, then $|f(x) - f(y)| < 1/n$. Now let $\{x_0, \dots, x_{m_n}\}$ be a partition of $[a, b]$ such that $|x_i - x_{i-1}| < \delta$ for each i . For $k = 1, 2, \dots, m_n - 1$, let $z_k^n \notin D$ and $|z_k^n - x_k| < \delta$. Then

$$|z_k^n - z_{k-1}^n| \leq |z_k^n - x_k| + |x_k - x_{k-1}| + |x_{k-1} - z_{k-1}^n| < 3\delta.$$

It follows that for each $x \in [a, b]$, $\left| \sum_{k=1}^{m_n} f(z_{k-1}^n) \mathcal{R}_{[z_{k-1}^n, z_k^n)}(x) - f(x) \right| < 1/n$. ■

Proposition 10.10.4 *Let f be a continuous function defined on \mathbb{R} . Also let F be an increasing function defined on \mathbb{R} . Then whenever c, d are not points of discontinuity of F and $[a, b] \supseteq [c, d]$, $\int_a^b f \mathcal{R}_{[c, d]} dF = \int f \mathcal{R}_{[c, d]} d\mu$. Here μ is the Lebesgue Stieltjes measure defined above.*

Proof: Since F is an increasing function it can have only countably many discontinuities. The reason for this is that the only kind of discontinuity it can have is where $F(x+) > F(x-)$. Now since F is increasing, the intervals $(F(x-), F(x+))$ for x a point of discontinuity are disjoint and so since each must contain a rational number and the rational numbers are countable, and therefore so are these intervals.

Let D denote this countable set of discontinuities of F . Then if $l, r \notin D, [l, r] \subseteq [a, b]$, it follows quickly from the definition of the Darboux Stieltjes integral that

$$\int_a^b \mathcal{X}_{[l,r]} dF = F(r) - F(l) = F(r-) - F(l-) = \mu([l, r]) = \int \mathcal{X}_{[l,r]} d\mu.$$

Now let $\{s_n\}$ be the sequence of step functions of Lemma 10.10.3 such that these step functions converge uniformly to f on $[c, d]$, say $\max_x |f(x) - s_n(x)| < 1/n$. Then

$$\left| \int (\mathcal{X}_{[c,d]} f - \mathcal{X}_{[c,d]} s_n) d\mu \right| \leq \int |\mathcal{X}_{[c,d]} (f - s_n)| d\mu \leq \frac{1}{n} \mu([c, d])$$

and $\left| \int_a^b (\mathcal{X}_{[c,d]} f - \mathcal{X}_{[c,d]} s_n) dF \right| \leq \int_a^b \mathcal{X}_{[c,d]} |f - s_n| dF < \frac{1}{n} (F(b) - F(a))$. Also if s_n is given by the formula of Lemma 10.10.3,

$$\begin{aligned} \int \mathcal{X}_{[c,d]} s_n d\mu &= \int \sum_{k=1}^{m_n} f(z_{k-1}^n) \mathcal{X}_{[z_{k-1}^n, z_k^n]} d\mu = \sum_{k=1}^{m_n} \int f(z_{k-1}^n) \mathcal{X}_{[z_{k-1}^n, z_k^n]} d\mu \\ &= \sum_{k=1}^{m_n} f(z_{k-1}^n) \mu([z_{k-1}^n, z_k^n]) = \sum_{k=1}^{m_n} f(z_{k-1}^n) (F(z_k^n) - F(z_{k-1}^n)) \\ &= \sum_{k=1}^{m_n} f(z_{k-1}^n) (F(z_k^n) - F(z_{k-1}^n)) = \sum_{k=1}^{m_n} \int_a^b f(z_{k-1}^n) \mathcal{X}_{[z_{k-1}^n, z_k^n]} dF = \int_a^b s_n dF. \end{aligned}$$

Therefore,

$$\begin{aligned} \left| \int \mathcal{X}_{[c,d]} f d\mu - \int_a^b \mathcal{X}_{[c,d]} f dF \right| &\leq \left| \int \mathcal{X}_{[c,d]} f d\mu - \int \mathcal{X}_{[c,d]} s_n d\mu \right| \\ &\quad + \left| \int \mathcal{X}_{[c,d]} s_n d\mu - \int_a^b s_n dF \right| + \left| \int_a^b s_n dF - \int_a^b \mathcal{X}_{[c,d]} f dF \right| \\ &\leq \frac{1}{n} \mu([c, d]) + \frac{1}{n} (F(b) - F(a)) \end{aligned}$$

and since n is arbitrary, this shows $\int f d\mu - \int_a^b f dF = 0$. ■

In particular, in the special case where F is continuous and f is continuous, $\int_a^b f dF = \int \mathcal{X}_{[a,b]} f d\mu$. Thus, if $F(x) = x$ so the Darboux Stieltjes integral is the usual integral from calculus, $\int_a^b f(t) dt = \int \mathcal{X}_{[a,b]} f d\mu$ where μ is the measure which comes from $F(x) = x$ as described above. This measure is often denoted by m . Thus when f is continuous $\int_a^b f(t) dt = \int \mathcal{X}_{[a,b]} f dm$ and so there is no problem in writing $\int_a^b f(t) dt$ for either the Lebesgue or the Riemann integral. Furthermore, when f is continuous, you can compute the Lebesgue integral by using the fundamental theorem of calculus because in this case, the two integrals are equal. ■

Note that as a special case, if f is continuous on \mathbb{R} and equals 0 off some finite interval I , written as $f \in C_c(\mathbb{R})$, then $\int_{\mathbb{R}} f dF = \int_{\mathbb{R}} f d\mu$ where the first integral is defined as $\int_c^d f dF$ where $(c, d) \supseteq I$. You could use the Riemann Stieltjes integral to define a positive linear functional on $C_c(\mathbb{R})$ as just explained and then it follows that the Lebesgue integral taken with respect to the Lebesgue Stieltjes measure above equals this functional on $C_c(\mathbb{R})$. This idea will be discussed more later in the abstract theory when such functionals will be shown to determine measures.

10.11 The Distribution Function

For $(\Omega, \mathcal{F}, \mu)$ a measure space, the integral of a nonnegative measurable function was defined earlier as $\int f d\mu \equiv \int_0^\infty \mu([f > t]) dt$. This idea will be developed more in this section.

Definition 10.11.1 Let $f \geq 0$ and suppose f is measurable. The distribution function is the function defined by $t \rightarrow \mu([f > t])$.

Lemma 10.11.2 If $\{f_n\}$ is an increasing sequence of functions converging pointwise to f then $\mu([f > t]) = \lim_{n \rightarrow \infty} \mu([f_n > t])$.

Proof: The sets, $[f_n > t]$ are increasing and their union is $[f > t]$ because if $f(\omega) > t$, then for all n large enough, $f_n(\omega) > t$ also. Therefore, the desired conclusion follows from properties of measures, the one which says that if $E_n \uparrow E$, then $\mu(E_n) \uparrow \mu(E)$. ■

Note how it was important to have strict inequality in the definition.

Lemma 10.11.3 Suppose $s \geq 0$ is a simple function, $s(\omega) \equiv \sum_{k=1}^n a_k \chi_{E_k}(\omega)$ where the a_k are the distinct nonzero values of s , $0 < a_1 < a_2 < \dots < a_n$ on the measurable sets E_k . Suppose ϕ is a C^1 function defined on $[0, \infty)$ which has the properties that $\phi(0) = 0$, and also that $\phi'(t) > 0$ for all t . Then

$$\int_0^\infty \phi'(t) \mu([s > t]) dm(t) = \int \phi(s) d\mu(s).$$

Proof: First note that if $\mu(E_k) = \infty$ for any k then both sides equal ∞ and so without loss of generality, assume $\mu(E_k) < \infty$ for all k . Letting $a_0 \equiv 0$, the left side equals

$$\begin{aligned} \sum_{k=1}^n \int_{a_{k-1}}^{a_k} \phi'(t) \mu([s > t]) dm(t) &= \sum_{k=1}^n \int_{a_{k-1}}^{a_k} \phi'(t) \sum_{i=k}^n \mu(E_i) dm \\ &= \sum_{k=1}^n \sum_{i=k}^n \mu(E_i) \int_{a_{k-1}}^{a_k} \phi'(t) dm = \sum_{k=1}^n \sum_{i=k}^n \mu(E_i) (\phi(a_k) - \phi(a_{k-1})) \\ &= \sum_{i=1}^n \mu(E_i) \sum_{k=1}^i (\phi(a_k) - \phi(a_{k-1})) = \sum_{i=1}^n \mu(E_i) \phi(a_i) = \int \phi(s) d\mu. \quad \blacksquare \end{aligned}$$

With this lemma the next theorem which is the main result follows easily.

Theorem 10.11.4 Let $f \geq 0$ be measurable and let ϕ be a C^1 function defined on $[0, \infty)$ which satisfies $\phi'(t) > 0$ for all $t > 0$ and $\phi(0) = 0$. Then

$$\int \phi(f) d\mu = \int_0^\infty \phi'(t) \mu([f > t]) dm.$$

Proof: By Theorem 9.1.6 on Page 239 there exists an increasing sequence of nonnegative simple functions, $\{s_n\}$ which converges pointwise to f . By the monotone convergence theorem and Lemma 10.11.2,

$$\begin{aligned} \int \phi(f) d\mu &= \lim_{n \rightarrow \infty} \int \phi(s_n) d\mu = \lim_{n \rightarrow \infty} \int_0^\infty \phi'(t) \mu([s_n > t]) dm \\ &= \int_0^\infty \phi'(t) \mu([f > t]) dm \blacksquare \end{aligned}$$

This theorem can be generalized to a situation in which ϕ is only increasing and continuous. In the generalization I will replace the symbol ϕ with F to coincide with earlier notation.

The following lemma and theorem say essentially that for F an increasing function equal to 0 at 0, $\int_{(0,\infty)} \mu([f > t]) dF = \int_\Omega F(f) d\mu$. I think it is particularly memorable if F is differentiable when it looks like what was just discussed. $\int_{(0,\infty)} \mu([f > t]) F'(t) dt = \int_\Omega F(f) d\mu$

Lemma 10.11.5 Suppose $s \geq 0$ is a simple function, $s(\omega) \equiv \sum_{k=1}^n a_k \chi_{E_k}(\omega)$ where the a_k are the distinct nonzero values of s , $a_1 < a_2 < \dots < a_n$. Suppose F is an increasing function defined on $[0, \infty)$, $F(0) = 0$, F being continuous at 0 from the right and continuous at every a_k . Then letting μ be a measure and $(\Omega, \mathcal{F}, \mu)$ a measure space,

$$\int_{(0,\infty)} \mu([s > t]) dv = \int_\Omega F(s) d\mu.$$

where the integral on the left is the Lebesgue integral for the Lebesgue Stieltjes measure ν which comes from the increasing function F as in Theorem 10.10.1 above.

Proof: This follows from the following computation. Since F is continuous at 0 and the values a_k ,

$$\begin{aligned} \int_0^\infty \mu([s > t]) dv(t) &= \sum_{k=1}^n \int_{(a_{k-1}, a_k]} \mu([s > t]) dv(t) \\ &= \sum_{k=1}^n \int_{(a_{k-1}, a_k]} \sum_{j=k}^n \mu(E_j) dF(t) = \sum_{j=1}^n \mu(E_j) \sum_{k=1}^j \nu((a_{k-1}, a_k]) \\ &= \sum_{j=1}^n \mu(E_j) \sum_{k=1}^j (F(a_k) - F(a_{k-1})) = \sum_{j=1}^n \mu(E_j) F(a_j) \equiv \int_\Omega F(s) d\mu \blacksquare \end{aligned}$$

Now here is the generalization to nonnegative measurable f .

Theorem 10.11.6 Let $f \geq 0$ be measurable with respect to \mathcal{F} , $(\Omega, \mathcal{F}, \mu)$ a measure space, and let F be an increasing continuous function defined on $[0, \infty)$ and $F(0) = 0$. Then $\int_\Omega F(f) d\mu = \int_{(0,\infty)} \mu([f > t]) dv(t)$ where ν is the Lebesgue Stieltjes measure determined by F as in Theorem 10.10.1 above.

Proof: By Theorem 9.1.6 on Page 239 there exists an increasing sequence of nonnegative simple functions, $\{s_n\}$ which converges pointwise to f . By the monotone convergence theorem and Lemma 10.11.5,

$$\int_\Omega F(f) d\mu = \lim_{n \rightarrow \infty} \int_\Omega F(s_n) d\mu = \lim_{n \rightarrow \infty} \int_{(0,\infty)} \mu([s_n > t]) dv = \int_{(0,\infty)} \mu([f > t]) dv \blacksquare$$

Note that the function $t \rightarrow \mu([f > t])$ is a decreasing function. Therefore, one can make sense of an improper Riemann Stieltjes integral $\int_0^\infty \mu([f > t]) dF(t)$. With more work, one can have this equal to the corresponding Lebesgue integral above.

10.12 Good Lambda Inequality

There is a very interesting and important inequality called the good lambda inequality (I am not sure if there is a bad lambda inequality.) which follows from the above theory of distribution functions. It involves the inequality

$$\mu([f > \beta\lambda] \cap [g \leq \delta\lambda]) \leq \phi(\delta) \mu([f > \lambda]) \quad (10.13)$$

for $\beta > 1$, nonnegative functions f, g and is supposed to hold for all small positive δ and $\phi(\delta) \rightarrow 0$ as $\delta \rightarrow 0$. Note the left side is small when g is large and f is small. The inequality involves dominating an integral involving f with one involving g as described below. As above, ν is the Lebesgue Stieltjes measure described above in terms of F , an increasing function. Is there any way to see the inequality in 10.13 might make sense? Look at the expression on the left. If δ is small enough, you might think that the intersection of the two sets would have smaller measure than $\mu([f > \lambda])$.

Theorem 10.12.1 *Let $(\Omega, \mathcal{F}, \mu)$ be a finite measure space and let F be a continuous increasing function defined on $[0, \infty)$ such that $F(0) = 0$. Suppose also that for every $\alpha > 1$, there exists a constant C_α such that for all $x \in [0, \infty)$, $F(\alpha x) \leq C_\alpha F(x)$. Also suppose f, g are nonnegative measurable functions and there exists $\beta > 1$, such that for all $\lambda > 0$ and $1 > \delta > 0$,*

$$\mu([f > \beta\lambda] \cap [g \leq \delta\lambda]) \leq \phi(\delta) \mu([f > \lambda]) \quad (10.14)$$

where $\lim_{\delta \rightarrow 0+} \phi(\delta) = 0$ and ϕ is increasing. Under these conditions, there exists a constant C depending only on β, ϕ such that

$$\int_{\Omega} F(f(\omega)) d\mu(\omega) \leq C \int_{\Omega} F(g(\omega)) d\mu(\omega).$$

Proof: Let $\beta > 1$ be as given above. First suppose f is bounded. This is so there can be no question of existence of the integrals. $\int_{\Omega} F(f) d\mu = \int_{\Omega} F\left(\beta \frac{f}{\beta}\right) d\mu \leq C_{\beta} \int_{\Omega} F\left(\frac{f}{\beta}\right) d\mu$. Let ν be the Lebesgue Stieltjes measure which comes from F , ($d\nu = dF$). From Theorem 10.11.6, $C_{\beta} \int_{\Omega} F\left(\frac{f}{\beta}\right) d\mu = C_{\beta} \int_0^\infty \mu([f/\beta > \lambda]) d\nu = C_{\beta} \int_0^\infty \mu([f > \beta\lambda]) d\nu$. Now using the given inequality,

$$\begin{aligned} \int_{\Omega} F(f) d\mu &= \\ & C_{\beta} \int_0^\infty \mu([f > \beta\lambda] \cap [g \leq \delta\lambda]) d\nu(\lambda) + C_{\beta} \int_0^\infty \mu([f > \beta\lambda] \cap [g > \delta\lambda]) d\nu(\lambda) \\ &\leq C_{\beta} \phi(\delta) \int_0^\infty \mu([f > \lambda]) d\nu(\lambda) + C_{\beta} \int_0^\infty \mu([g > \delta\lambda]) d\nu(\lambda) \\ &\leq C_{\beta} \phi(\delta) \int_{\Omega} F(f) d\mu + C_{\beta} \int_{\Omega} F\left(\frac{g}{\delta}\right) d\mu \end{aligned}$$

Now choose δ small enough that $C_\beta \phi(\delta) < \frac{1}{2}$ and then subtract the first term on the right in the above from both sides. It follows from the properties of F again that

$$\int_{\Omega} F(f) d\mu \leq 2C_\beta C_{\delta^{-1}} \int_{\Omega} F(g) d\mu. \quad (10.15)$$

This establishes the inequality in the case where f is bounded.

In general, let $f_n = \min(f, n)$. For $n \leq \lambda$, the inequality

$$\mu([f > \beta\lambda] \cap [g \leq \delta\lambda]) \leq \phi(\delta) \mu([f > \lambda])$$

holds with f replaced with f_n because both sides equal 0 thanks to $\beta > 1$. If $n > \lambda$, then $[f > \lambda] = [f_n > \lambda]$ and so the inequality still holds because in this case,

$$\begin{aligned} \mu([f_n > \beta\lambda] \cap [g \leq \delta\lambda]) &\leq \mu([f > \beta\lambda] \cap [g \leq \delta\lambda]) \\ &\leq \phi(\delta) \mu([f > \lambda]) = \phi(\delta) \mu([f_n > \lambda]) \end{aligned}$$

Therefore, 10.14 is valid with f replaced with f_n . Now pass to the limit in $\int_{\Omega} F(f_n) d\mu \leq 2C_\beta C_{\delta^{-1}} \int_{\Omega} F(g_n) d\mu$ as $n \rightarrow \infty$ and use the monotone convergence theorem. ■

10.13 Radon Nikodym Theorem

Let μ, ν be two finite measures on the measurable space (Ω, \mathcal{F}) and let $\alpha \geq 0$. Let $\lambda \equiv \nu - \alpha\mu$. Then it is clear that if $\{E_i\}_{i=1}^{\infty}$ are disjoint sets of \mathcal{F} , then $\lambda(\cup_i E_i) = \sum_{i=1}^{\infty} \lambda(E_i)$ and that the series converges. The next proposition is fairly obvious.

Proposition 10.13.1 *Let $(\Omega, \mathcal{F}, \lambda)$ be a measure space and let $\lambda : \mathcal{F} \rightarrow [0, \infty)$ be a measure. Then λ is a finite measure.*

Proof: Since $\lambda(\Omega) < \infty$ this is a finite measure. ■

Definition 10.13.2 *Let (Ω, \mathcal{F}) be a measurable space and let $\lambda : \mathcal{F} \rightarrow \mathbb{R}$ satisfy: If $\{E_i\}_{i=1}^{\infty}$ are disjoint sets of \mathcal{F} , then $\lambda(\cup_i E_i) = \sum_{i=1}^{\infty} \lambda(E_i)$ and the series converges. Such a real valued function is called a signed measure. In this context, a set $E \in \mathcal{F}$ is called positive if whenever F is a measurable subset of E , it follows $\lambda(F) \geq 0$. A negative set is defined similarly. Note that this requires $\lambda(\Omega) \in \mathbb{R}$.*

Lemma 10.13.3 *The countable union of disjoint positive sets is positive.*

Proof: Let E_i be positive and consider $E \equiv \cup_{i=1}^{\infty} E_i$. If $A \subseteq E$ with A measurable, then $A \cap E_i \subseteq E_i$ and so $\lambda(A \cap E_i) \geq 0$. Hence $\lambda(A) = \sum_i \lambda(A \cap E_i) \geq 0$. ■

Lemma 10.13.4 *Let λ be a signed measure on (Ω, \mathcal{F}) . If $E \in \mathcal{F}$ with $0 < \lambda(E)$, then E has a measurable subset which is positive.*

Proof: If every measurable subset F of E has $\lambda(F) \geq 0$, then E is positive and we are done. Otherwise there exists measurable $F \subseteq E$ with $\lambda(F) < 0$. Let the elements of \mathfrak{F} consist of sets of disjoint sets of measurable subsets of E each of which has measure less than 0. Partially order \mathfrak{F} by set inclusion. By the Hausdorff maximal theorem, Theorem 2.8.4, there is a maximal chain \mathcal{C} . Then $\cup \mathcal{C}$ is a set consisting of disjoint measurable sets $F \in \mathcal{F}$ such that $\lambda(F) < 0$. Since each set in $\cup \mathcal{C}$ has measure strictly less than 0, it follows

that $\cup \mathcal{C}$ is a countable set, $\{F_i\}_{i=1}^\infty$. Otherwise, there would exist an infinite subset of $\cup \mathcal{C}$ with each set having measure less than $-\frac{1}{n}$ for some $n \in \mathbb{N}$ so λ would not be real valued. Letting $F = \cup_i F_i$, then $E \setminus F$ has no measurable subsets S for which $\lambda(S) < 0$ since, if it did, \mathcal{C} would not have been maximal. Thus $E \setminus F$ is positive. ■

A major result is the following, called a Hahn decomposition.

Theorem 10.13.5 *Let λ be a signed measure on a measurable space (Ω, \mathcal{F}) . Then there are disjoint measurable sets P, N such that P is a positive set, N is a negative set, and $P \cup N = \Omega$.*

Proof: If Ω is either positive or negative, there is nothing to show, so suppose Ω is neither positive nor negative. \mathfrak{F} will consist of collections of disjoint measurable sets F such that $\lambda(F) > 0$. Thus each element of \mathfrak{F} is necessarily countable. Partially order \mathfrak{F} by set inclusion and use the Hausdorff maximal theorem to get \mathcal{C} a maximal chain. Then, as in the above lemma, $\cup \mathcal{C}$ is countable, say $\{P_i\}_{i=1}^\infty$ because $\lambda(F) > 0$ for each $F \in \cup \mathcal{C}$ and λ has values in \mathbb{R} . The sets in $\cup \mathcal{C}$ are disjoint because if A, B are two of them, then they are both in a single element of \mathcal{C} . Letting $P \equiv \cup_i P_i$, and $N = P^C$, it follows from Lemma 10.13.3 that P is positive. It is also the case that N must be negative because otherwise, \mathcal{C} would not be maximal. ■

Clearly a Hahn decomposition is not unique. For example, you could have obtained a different Hahn decomposition if you had considered disjoint negative sets F for which $\lambda(F) < 0$ in the above argument.

Let $k \in \mathbb{N}$, $\{\alpha_n^k\}_{n=0}^\infty$ be equally spaced points $\alpha_n^k = 2^{-k}n$. Then $\alpha_{2n}^k = 2^{-k}(2n) = 2^{-(k-1)}n \equiv \alpha_n^{k-1}$ and $\alpha_{2n}^{k+1} \equiv 2^{-(k+1)}2n = \alpha_n^k$. Similarly $N_{2n}^{k+1} = N_n^k$ because these depend on the α_n^k . Also let (P_n^k, N_n^k) be a Hahn decomposition for the signed measure $\nu - \alpha_n^k \mu$ where ν, μ are two finite measures. Now from the definition, $N_{n+1}^k \setminus N_n^k = N_{n+1}^k \cap P_n^k$. Also, $N_n \subseteq N_{n+1}$ for each n and we can take $N_0 = \emptyset$. Then $\{N_{n+1}^k \setminus N_n^k\}_{n=0}^\infty$ covers all of Ω except possibly for a set of μ measure 0.

Lemma 10.13.6 *Let $S \equiv \Omega \setminus (\cup_n N_n^k) = \Omega \setminus (\cup_n N_n^l)$ for any l . Then $\mu(S) = 0$.*

Proof: $S = \cap_n P_n^k$ so for all n , $\nu(S) - \alpha_n^k \mu(S) \geq 0$. But letting $n \rightarrow \infty$, it must be that $\mu(S) = 0$. ■

As just noted, if $E \subseteq N_{n+1}^k \setminus N_n^k$, then

$$\nu(E) - \alpha_n^k \mu(E) \geq 0 \geq \nu(E) - \alpha_{n+1}^k \mu(E), \text{ so } \alpha_{n+1}^k \mu(E) \geq \nu(E) \geq \alpha_n^k \mu(E) \quad (10.16)$$

$$\begin{array}{c} \boxed{\begin{array}{c} N_{n+1}^k \\ \boxed{N_n^k} \\ \alpha_{n+1}^k \mu(E) \geq \nu(E) \geq \alpha_n^k \mu(E) \end{array}} \end{array}$$

Then define $f^k(\omega) \equiv \sum_{n=0}^\infty \alpha_n^k \mathcal{X}_{\Delta_n^k}(\omega)$ where $\Delta_m^k \equiv N_{m+1}^k \setminus N_m^k$. Thus,

$$\begin{aligned} f^k &= \sum_{n=0}^\infty \alpha_{2n}^{k+1} \mathcal{X}_{(N_{2n+2}^{k+1} \setminus N_{2n}^{k+1})} = \sum_{n=0}^\infty \alpha_{2n}^{k+1} \mathcal{X}_{\Delta_{2n+1}^{k+1}} + \sum_{n=0}^\infty \alpha_{2n}^{k+1} \mathcal{X}_{\Delta_{2n}^{k+1}} \\ &\leq \sum_{n=0}^\infty \alpha_{2n+1}^{k+1} \mathcal{X}_{\Delta_{2n+1}^{k+1}} + \sum_{n=0}^\infty \alpha_{2n}^{k+1} \mathcal{X}_{\Delta_{2n}^{k+1}} = f^{k+1} \end{aligned} \quad (10.17)$$

Thus $k \rightarrow f^k(\omega)$ is increasing. Let $f(\omega) \equiv \lim_{k \rightarrow \infty} f^k(\omega)$. Thus $f = 0$ on S . Now let E be measurable. Thus $\mu(E) = \mu(E \cap S) + \mu(E \cap S^C)$, similar for λ and $\lambda(E \cap S^C) = \sum_n \lambda(E \cap S^C \cap \Delta_n^k)$. To save space, let $\tilde{E} \equiv E \cap S^C$. Then using 10.16

$$\begin{aligned} \int \mathcal{X}_{\tilde{E}} f^k d\mu &\leq \sum_{n=0}^{\infty} \alpha_{n+1}^k \mu(\tilde{E} \cap \Delta_n^k) \leq \sum_{n=0}^{\infty} \alpha_n^k \mu(\tilde{E} \cap \Delta_n^k) + \sum_{n=0}^{\infty} 2^{-k} \mu(\tilde{E} \cap \Delta_n^k) \\ &\leq \sum_{n=0}^{\infty} \nu(\tilde{E} \cap \Delta_n^k) + 2^{-k} \mu(\tilde{E}) = \nu(\tilde{E}) + 2^{-k} \mu(\tilde{E}) \leq \int \mathcal{X}_{\tilde{E}} f^k d\mu + 2^{-k} \mu(\tilde{E}) \quad (10.18) \end{aligned}$$

From the monotone convergence theorem it follows $\nu(\tilde{E}) = \int \mathcal{X}_{\tilde{E}} f d\mu = \int \mathcal{X}_E f d\mu$.

This proves most of the following theorem which is the Radon Nikodym theorem.

Theorem 10.13.7 *Let ν and μ be finite measures defined on a measurable space (Ω, \mathcal{F}) . Then there exists a measurable set S with $\mu(S) = 0$ and a nonnegative measurable function $\omega \rightarrow f(\omega)$ such that $\nu(E) = \int_E f d\mu + \nu(E \cap S)$.*

Proof: Let S be defined in Lemma 10.13.6 so $S \equiv \Omega \setminus (\cup_n N_n^k)$ and $\mu(S) = 0$. If $E \in \mathcal{F}$, and f as described above,

$$\nu(E) = \nu(E \cap S^C) + \nu(E \cap S) = \int_{E \cap S^C} f d\mu + \nu(E \cap S) = \int_E f d\mu + \nu(E \cap S)$$

Thus if $E \subseteq S^C$, we have $\nu(E) = \int_E f d\mu$. ■

Definition 10.13.8 *Let μ, ν be finite measures on (Ω, \mathcal{F}) . Then $\nu \ll \mu$ means that whenever $\mu(E) = 0$, it follows that $\nu(E) = 0$.*

Sometimes people write $f = \frac{d\nu}{d\mu}$, in the case $\nu \ll \mu$ and this is called the Radon Nikodym derivative.

Proposition 10.13.9 *If ν, μ are finite measures and $\nu \ll \mu$, then there exists nonnegative measurable f such that $\nu(E) = \int_E f d\mu$.*

Proof: In Theorem 10.13.7, $\nu(E \cap S) = 0$ because $\mu(S) = 0$ and $\nu \ll \mu$. ■

Definition 10.13.10 *Let S be in the above theorem. Then*

$$\nu_{||}(E) \equiv \nu(E \cap S^C) = \int_{E \cap S^C} f d\mu = \int_E f d\mu$$

while $\nu_{\perp}(E) \equiv \nu(E \cap S)$. Thus $\nu_{||} \ll \mu$ and ν_{\perp} is nonzero only on sets which are contained in S which has μ measure 0.

Corollary 10.13.11 *In the above situation, let λ be a signed measure and let $\lambda \ll \mu$ meaning that if $\mu(E) = 0 \Rightarrow \lambda(E) = 0$. Here assume that μ is a finite measure. Then there exists $h \in L^1$ such that $\lambda(E) = \int_E h d\mu$.*

Proof: Let $P \cup N$ be a Hahn decomposition of λ . Let

$$\lambda_+(E) \equiv \lambda(E \cap P), \quad \lambda_-(E) \equiv -\lambda(E \cap N).$$

Then both λ_+ and λ_- are absolutely continuous measures and so there are nonnegative h_+ and h_- with $\lambda_-(E) = \int_E h_- d\mu$ and a similar equation for λ_+ . Then $0 \leq -\lambda(\Omega \cap N) \leq \lambda_-(\Omega) < \infty$, similar for λ_+ so both of these measures are necessarily finite. Hence both h_- and h_+ are in L^1 so $h \equiv h_+ - h_-$ is also in L^1 and $\lambda(E) = \lambda_+(E) - \lambda_-(E) = \int_E (h_+ - h_-) d\mu$. ■

Proposition 10.13.12 *This Lebesgue decomposition is unique. If f, \hat{f} both work in Theorem 10.13.7, then $f = \hat{f}$ μ a.e. This function $f \in L^1(\Omega)$, $\int_{\Omega} f d\mu < \infty$.*

Proof: Say $\nu_{\parallel} + \nu_{\perp} = \hat{\nu}_{\parallel} + \hat{\nu}_{\perp}$. Then $\nu_{\parallel}(E) - \hat{\nu}_{\parallel}(E) = \nu_{\perp}(E) - \hat{\nu}_{\perp}(E)$. If $\mu(E) = 0$, then the left side is also 0 and so $\nu_{\perp}(E) - \hat{\nu}_{\perp}(E) = 0$. But then for any E ,

$$\nu_{\perp}(E) - \hat{\nu}_{\perp}(E) = \int_E h d\mu \quad (10.19)$$

for h a function in $L^1(\Omega)$. This is because $\nu_{\perp} - \hat{\nu}_{\perp}$ is a signed measure $\lambda \ll \mu$ and Corollary 10.13.11. From the above, if S, \hat{S} are the exceptional sets of μ measure zero,

$$\nu_{\perp}(E) - \hat{\nu}_{\perp}(E) = \nu_{\perp}(E \cup (S \cup \hat{S})) - \hat{\nu}_{\perp}(E \cup (S \cup \hat{S})) = \int_{E \cup (S \cup \hat{S})} h d\mu = 0 \quad (10.20)$$

because $\mu(S \cup \hat{S}) = 0$ and so the right side of 10.19 must be 0 after all, so $\nu_{\perp}(E) = \hat{\nu}_{\perp}(E)$. It follows that $\nu_{\parallel} = \hat{\nu}_{\parallel}$ also. Now in Theorem 10.13.7, if you have two f, \hat{f} which work, then

$$\nu_{\parallel}(E) = \int_E f d\mu = \int_E \hat{f} d\mu = \hat{\nu}_{\parallel}(E) \quad (10.21)$$

and so, $f = \hat{f}$ a.e. because you can apply this equation to $E_n \equiv [f - \hat{f} > 1/n]$ and conclude that

$$0 = \int_{E_n} f - \hat{f} d\mu \geq \frac{1}{n} \mu(E_n) = 0 \quad (10.22)$$

so $\mu([f - \hat{f} > 0]) = \cup_m \mu(E_m) = 0$. Similarly $\mu([\hat{f} - f > 0]) = 0$. ■

This unique decomposition of a measure ν into the sum of two measures, one absolutely continuous with respect to μ and the other supported on a set of μ measure zero is called the Lebesgue decomposition.

Definition 10.13.13 *A measure space $(\Omega, \mathcal{F}, \mu)$ is σ finite if there are countably many measurable sets $\{\Omega_n\}$ such that μ is finite on measurable subsets of Ω_n .*

There is a routine corollary of the above theorem.

Corollary 10.13.14 *Suppose μ, ν are both σ finite measures defined on (Ω, \mathcal{F}) . Then a similar conclusion to the above theorem can be obtained.*

$$\nu(E) = \int_E f d\mu + \nu(E \cap S), \quad \mu(S) = 0 \quad (10.23)$$

for f a nonnegative measurable function. If $\nu(\Omega) < \infty$, then $f \in L^1(\Omega)$. This f is unique up to a set of μ measure zero.

Proof: Since both μ, ν are σ finite, there are $\{\tilde{\Omega}_k\}_{k=1}^{\infty}$ such that $\nu(\tilde{\Omega}_k), \mu(\tilde{\Omega}_k)$ are finite. Let $\Omega_0 = \emptyset$ and $\Omega_k \equiv \tilde{\Omega}_k \setminus \left(\bigcup_{j=0}^{k-1} \tilde{\Omega}_j\right)$ so that μ, ν are finite on Ω_k and the Ω_k are disjoint. Let \mathcal{F}_k be the measurable subsets of Ω_k , equivalently the intersections with Ω_k with sets of \mathcal{F} . Now let $\nu_k(E) \equiv \nu(E \cap \Omega_k)$, similar for μ_k . By Theorem 10.13.7, there exists $S_k \subseteq \Omega_k$, and f_k as described there, unique up to sets of μ measure 0. Thus $\mu_k(S_k) = 0$ and $\nu_k(E) = \int_{E \cap \Omega_k} f_k d\mu_k + \nu_k(E \cap S_k)$. Now let $f(\omega) \equiv f_k(\omega)$ for $\omega \in \Omega_k$. Thus

$$\nu(E \cap \Omega_k) = \nu(E \cap (\Omega_k \setminus S_k)) + \nu(E \cap S_k) = \int_{E \cap \Omega_k} f d\mu + \nu(E \cap S_k) \quad (10.24)$$

Summing over all k , $\nu(E) = \nu(E \cap S^C) + \nu(E \cap S) = \int_E f d\mu + \nu(E \cap S)$. In particular, if $\nu \ll \mu$, then $\nu(E \cap S) = 0$ and so $\nu(E) = \int_E f d\mu$. The last claim is obvious from 10.23. ■

10.14 Iterated Integrals

This is about what can be said for the σ algebra of product measurable sets. First it is necessary to define what this means.

Definition 10.14.1 A measure space $(\Omega, \mathcal{F}, \mu)$ is called σ finite if there are measurable subsets Ω_n such that $\mu(\Omega_n) < \infty$ and $\Omega = \bigcup_{n=1}^{\infty} \Omega_n$.

Next is a σ algebra which comes from two σ algebras.

Definition 10.14.2 Let $(X, \mathcal{E}), (Y, \mathcal{F})$ be measurable spaces. That is, a set with a σ algebra of subsets. Then $\mathcal{E} \times \mathcal{F}$ will be the smallest σ algebra which contains the measurable rectangles, sets of the form $E \times F$ where $E \in \mathcal{E}, F \in \mathcal{F}$. The sets in this new σ algebra are called product measurable sets.

Definition 10.14.3 Given two finite measure spaces, (X, \mathcal{E}, μ) and (Y, \mathcal{F}, ν) , one can define a new measure $\mu \times \nu$ defined on $\mathcal{E} \times \mathcal{F}$ by specifying what it does to measurable rectangles as follows:

$$(\mu \times \nu)(A \times B) = \mu(A) \nu(B)$$

whenever $A \in \mathcal{E}$ and $B \in \mathcal{F}$.

We also have the following important proposition which holds in every context independent of any measure.

Proposition 10.14.4 Let $E \subseteq \mathcal{E} \times \mathcal{F}$ be product measurable $\mathcal{E} \times \mathcal{F}$ where \mathcal{E} is a σ algebra of sets of X and \mathcal{F} is a σ algebra of sets of Y . then if $E_x \equiv \{y \in Y : (x, y) \in E\}$ and $E_y \equiv \{x \in X : (x, y) \in E\}$, then $E_x \in \mathcal{F}$ and $E_y \in \mathcal{E}$.

Proof: It is obvious that if \mathcal{H} is the measurable rectangles, then the conclusion of the proposition holds. If \mathcal{G} consists of the sets of $\mathcal{E} \times \mathcal{F}$ for which the proposition holds, then it is clearly closed with respect to countable disjoint unions and complements. This is obvious in the case of a countable disjoint union since $(\bigcup_i E^i)_x = \bigcup_i E^i_x$, similar for y . As to complement, if $E \in \mathcal{G}$, then $E_x \in \mathcal{F}$ and so $(E^c)_x = (E_x)^c \in \mathcal{F}$. It is similar for y . By Dynkin's lemma, $\mathcal{G} \supseteq \mathcal{E} \times \mathcal{F}$. However \mathcal{G} was defined as a subset of $\mathcal{E} \times \mathcal{F}$ so these are equal. ■

Let (X, \mathcal{E}, μ) and (Y, \mathcal{F}, ν) be two finite measure spaces. Define \mathcal{H} to be the set of measurable rectangles, $A \times B, A \in \mathcal{E}$ and $B \in \mathcal{F}$. Let

$$\mathcal{G} \equiv \left\{ E \subseteq X \times Y : \int_Y \int_X \mathcal{X}_E d\mu d\nu = \int_X \int_Y \mathcal{X}_E d\nu d\mu \right\} \quad (10.25)$$

where in the above, part of the requirement is for all integrals to make sense.

Then $\mathcal{H} \subseteq \mathcal{G}$. This is obvious.

Next I want to show that if $E \in \mathcal{G}$ then $E^c \in \mathcal{G}$. Observe $\mathcal{X}_{E^c} = 1 - \mathcal{X}_E$ and so

$$\begin{aligned} \int_Y \int_X \mathcal{X}_{E^c} d\mu d\nu &= \int_Y \int_X (1 - \mathcal{X}_E) d\mu d\nu = \int_X \int_Y (1 - \mathcal{X}_E) d\nu d\mu \\ &= \int_X \int_Y \mathcal{X}_{E^c} d\nu d\mu \end{aligned}$$

which shows that if $E \in \mathcal{G}$, then $E^c \in \mathcal{G}$.

Next I want to show \mathcal{G} is closed under countable unions of disjoint sets of \mathcal{G} . Let $\{A_i\}$ be a sequence of disjoint sets from \mathcal{G} . Then, using the monotone convergence theorem as needed,

$$\begin{aligned} \int_Y \int_X \mathcal{X}_{\bigcup_{i=1}^{\infty} A_i} d\mu dv &= \int_Y \int_X \sum_{i=1}^{\infty} \mathcal{X}_{A_i} d\mu dv = \int_Y \sum_{i=1}^{\infty} \int_X \mathcal{X}_{A_i} d\mu dv \\ &= \sum_{i=1}^{\infty} \int_Y \int_X \mathcal{X}_{A_i} d\mu dv = \sum_{i=1}^{\infty} \int_X \int_Y \mathcal{X}_{A_i} dv d\mu \\ &= \int_X \sum_{i=1}^{\infty} \int_Y \mathcal{X}_{A_i} dv d\mu = \int_X \int_Y \sum_{i=1}^{\infty} \mathcal{X}_{A_i} dv d\mu = \int_X \int_Y \mathcal{X}_{\bigcup_{i=1}^{\infty} A_i} dv d\mu, \end{aligned} \quad (10.26)$$

Thus \mathcal{G} is closed with respect to countable disjoint unions.

From Lemma 9.3.2, $\mathcal{G} \supseteq \sigma(\mathcal{K})$, the smallest σ algebra containing \mathcal{K} . Also the computation in 10.26 implies that on $\sigma(\mathcal{K})$ one can define a measure, denoted by $\mu \times \nu$ and that for every $E \in \sigma(\mathcal{K})$,

$$(\mu \times \nu)(E) = \int_Y \int_X \mathcal{X}_E d\mu dv = \int_X \int_Y \mathcal{X}_E dv d\mu. \quad (10.27)$$

with each iterated integral making sense.

Next is product measure. First is the case of finite measures. Then this will extend to σ finite measures. The following theorem is Fubini's theorem.

Theorem 10.14.5 *Let $f : X \times Y \rightarrow [0, \infty]$ be measurable with respect to the σ algebra, $\sigma(\mathcal{K}) \equiv \mathcal{E} \times \mathcal{F}$ just defined and let $\mu \times \nu$ be the product measure of 10.27 where μ and ν are finite measures on (X, \mathcal{E}) and (Y, \mathcal{F}) respectively. Then*

$$\int_{X \times Y} f d(\mu \times \nu) = \int_Y \int_X f d\mu dv = \int_X \int_Y f dv d\mu.$$

Proof: Let $\{s_n\}$ be an increasing sequence of $\sigma(\mathcal{K}) \equiv \mathcal{E} \times \mathcal{F}$ measurable simple functions which converges pointwise to f . The above equation holds for s_n in place of f from what was shown above. The final result follows from passing to the limit and using the monotone convergence theorem. ■

Of course one can generalize right away to measures which are only σ finite. This is also called Fubini's theorem.

Definition 10.14.6 *Let $(X, \mathcal{E}, \mu), (Y, \mathcal{F}, \nu)$ both be σ finite. Thus there exist disjoint measurable X_n with $\bigcup_{n=1}^{\infty} X_n = X$ and disjoint measurable Y_n with $\bigcup_{n=1}^{\infty} Y_n = Y$ such that μ, ν restricted to X_n, Y_n respectively are finite measures. Let \mathcal{E}_n be intersections of sets of \mathcal{E} with X_n and \mathcal{F}_n similarly defined. Then letting \mathcal{K} consist of all measurable rectangles $A \times B$ for $A \in \mathcal{E}, B \in \mathcal{F}$, and letting $\mathcal{E} \times \mathcal{F} \equiv \sigma(\mathcal{K})$ define the product measure of E contained in this σ algebra as $(\mu \times \nu)(E) \equiv \sum_n \sum_m (\mu_n \times \nu_m)(E \cap (X_n \times Y_m))$.*

Lemma 10.14.7 *The above definition yields a well defined measure on $\mathcal{E} \times \mathcal{F}$.*

Proof: This follows from the standard theorems about sums of nonnegative numbers. See Theorem 2.5.4. For example if you have two other disjoint sequences X_k, Y_l on which the measures are finite, then

$$\begin{aligned} (\mu \times \nu)(E) &= \sum_n \sum_m \sum_k \sum_l (\mu_n \times \nu_m)(E \cap (X_n \cap X_k \times Y_m \cap Y_l)) \\ &= \sum_k \sum_l \sum_n \sum_m (\mu_k \times \nu_l)(E \cap (X_n \cap X_k \times Y_m \cap Y_l)) \end{aligned}$$

and so the definition with respect to the two different increasing sequences gives the same thing. Thus the definition is well defined. $(\mu \times \nu)$ is a measure because if the E_i are disjoint $\mathcal{E} \times \mathcal{F}$ measurable sets and $E = \cup_i E_i$,

$$\begin{aligned} (\mu \times \nu)(E) &\equiv \sum_n \sum_m (\mu_n \times \nu_m)(\cup_i E_i \cap (X_n \times Y_m)) = \sum_n \sum_m \sum_i (\mu_n \times \nu_m)(E_i \cap (X_n \times Y_m)) \\ &= \sum_i \sum_n \sum_m (\mu_n \times \nu_m)(E_i \cap (X_n \times Y_m)) \equiv \sum_i (\mu \times \nu)(E_i) \blacksquare \end{aligned}$$

Theorem 10.14.8 *Let $f : X \times Y \rightarrow [0, \infty]$ be measurable with respect to the σ algebra, $\sigma(\mathcal{K})$ just defined as the smallest σ algebra containing the measurable rectangles, and let $\mu \times \nu$ be the product measure of 10.27 where μ and ν are σ finite measures on (X, \mathcal{E}) and (Y, \mathcal{F}) respectively. (10.14.1) Then*

$$\int_{X \times Y} f d(\mu \times \nu) = \int_Y \int_X f d\mu d\nu = \int_X \int_Y f d\nu d\mu. \quad (10.28)$$

Proof: Letting $E \in \mathcal{E} \times \mathcal{F}$,

$$\begin{aligned} \int_{X \times Y} \mathcal{X}_E d(\mu \times \nu) &\equiv (\mu \times \nu)(E) \equiv \sum_n \sum_m (\mu_n \times \nu_m)(E \cap (X_n \times Y_m)) \\ &= \sum_n \sum_m \int_{Y_n} \int_{X_n} \mathcal{X}_E d\mu_n d\nu_m = \int_Y \int_X \mathcal{X}_E d\mu d\nu \end{aligned}$$

the last coming from a use of the monotone convergence theorem applied to sums. It follows that 10.28 holds for simple functions and then from monotone convergence theorem and Theorem 9.1.6, it holds for nonnegative $\mathcal{E} \times \mathcal{F}$ measurable functions. ■

It is also useful to note that all the above holds for $\prod_{i=1}^p X_i$ in place of $X \times Y$ and μ_i a measure on \mathcal{E}_i a σ algebra of sets of X_i . You would simply modify the definition of \mathcal{G} in 10.25 including all permutations for the iterated integrals and for \mathcal{K} you would use sets of the form $\prod_{i=1}^p A_i$ where A_i is measurable. Everything goes through exactly as above.

Thus the following is mostly obtained.

Theorem 10.14.9 *Let $\{(X_i, \mathcal{E}_i, \mu_i)\}_{i=1}^p$ be σ finite measure spaces and $\prod_{i=1}^p \mathcal{E}_i$ denotes the smallest σ algebra which contains the measurable boxes of the form $\prod_{i=1}^p A_i$ where $A_i \in \mathcal{E}_i$. Then there exists a measure λ defined on $\prod_{i=1}^p \mathcal{E}_i$ such that if $f : \prod_{i=1}^p X_i \rightarrow [0, \infty]$ is $\prod_{i=1}^p \mathcal{E}_i$ measurable, (i_1, \dots, i_p) is any permutation of $(1, \dots, p)$, then*

$$\int f d\lambda = \int_{X_{i_p}} \cdots \int_{X_{i_1}} f d\mu_{i_1} \cdots d\mu_{i_p} \quad (10.29)$$

If each X_i is a complete separable metric space such that μ_i is finite on balls and \mathcal{E}_i contains $\mathcal{B}(X_i)$, the Borel sets of X_i , then λ is a regular measure on a σ algebra of sets of $\prod_{i=1}^p X_i$ with the metric given by $d(\mathbf{x}, \mathbf{y}) \equiv \max\{d(x_i, y_i) : x_i, y_i \in X_i\}$, which includes the Borel sets.

Proof: It remains to verify the last claim. This is because all sets $\prod_{i=1}^p B(\xi_i, r)$ are contained in $\prod_{i=1}^p \mathcal{E}_i$ and are the open balls for the topology of $\prod_{i=1}^p X_i$. Then by separability of each X_i , the product $\prod_{i=1}^p X_i$ is also separable and so this product with the above metric is completely separable. Thus every open set is the countable union of these sets so open sets are in $\prod_{i=1}^p \mathcal{E}_i$ which consequently contains the Borel sets. Now from Corollary 9.8.9, λ is regular because it is finite on balls. ■

The conclusion 10.29 is called Fubini's theorem. More generally

Theorem 10.14.10 Suppose, in the situation of Theorem 10.14.9 $f \in L^1$ with respect to the measure λ . Then 10.29 continues to hold.

Proof: It suffices to prove this for f having real values because if this is shown the general case is obtained by taking real and imaginary parts. Since $f \in L^1(\prod_{i=1}^p X_i)$, $\int |f| d\lambda < \infty$ and so both $\frac{1}{2}(|f| + f)$ and $\frac{1}{2}(|f| - f)$ are in $L^1(\prod_{i=1}^p X_i)$ and are each non-negative. Hence from Theorem 10.14.9,

$$\begin{aligned} \int f d\lambda &= \int \left[\frac{1}{2}(|f| + f) - \frac{1}{2}(|f| - f) \right] d\lambda = \int \frac{1}{2}(|f| + f) d\lambda - \int \frac{1}{2}(|f| - f) d\lambda \\ &= \int \cdots \int \frac{1}{2}(|f| + f) d\mu_{i_1} \cdots d\mu_{i_p} - \int \cdots \int \frac{1}{2}(|f| - f) d\mu_{i_1} \cdots d\mu_{i_p} \\ &= \int \cdots \int \frac{1}{2}(|f| + f) - \frac{1}{2}(|f| - f) d\mu_{i_1} \cdots d\mu_{i_p} = \int \cdots \int f d\mu_{i_1} \cdots d\mu_{i_p} \quad \blacksquare \end{aligned}$$

The following corollary is a convenient way to verify the hypotheses of the above theorem.

Corollary 10.14.11 Suppose f is measurable with respect to $\prod_{i=1}^p \mathcal{E}_i$ and suppose for some permutation, (i_1, \dots, i_p) , $\int \cdots \int |f| d\mu_{i_1} \cdots d\mu_{i_p} < \infty$. Then $f \in L^1(\prod_{i=1}^p X_i)$.

Proof: By Theorem 10.14.9, $\int_{\mathbb{R}^p} |f| d\lambda = \int \cdots \int |f(x)| d\mu_{i_1} \cdots d\mu_{i_p} < \infty$ and so f is in $L^1(\mathbb{R}^p)$. \blacksquare

You can of course consider the completion of a product measure by using the outer measure approach described earlier. This could have been used to get Lebesgue measure.

If you have $f \geq 0$ and you consider the completion of $(\prod_{i=1}^p X_i, \prod_{i=1}^p \mathcal{E}_i, \prod_{i=1}^p \mu_i)$ denoted by $(\prod_{i=1}^p X_i, \overline{\prod_{i=1}^p \mathcal{E}_i}, \overline{\prod_{i=1}^p \mu_i})$ and f is $\overline{\prod_{i=1}^p \mathcal{E}_i}$ measurable, then the procedure for completing a measure space implies there are $h \geq f \geq g$ where h, g are both $\prod_{i=1}^p \mathcal{E}_i$ measurable and $f = g = h$ a.e. relative to the complete measure. Here $(X_i, \mathcal{E}_i, \mu_i)$ is finite or σ finite. Then you can define the iterated integral of f as that single number which is between the iterated integrals of g and h both of which make sense. Thus one can make sense of an iterated integral even if the function is not product measurable as long as the function is measurable in the completion of the product measure space. See Problem 18 on Page 354.

Given a finite measure μ defined on the product measurable sets, there exists a decomposition into iterated integrals as described in the following theorem.

Theorem 10.14.12 Let (X, \mathcal{E}) , and (Y, \mathcal{F}) be measurable spaces and let $\mathcal{E} \times \mathcal{F} \equiv \sigma(\mathcal{H})$ where \mathcal{H} denotes the measurable rectangles $E \times F$ for $E \in \mathcal{E}$ and $F \in \mathcal{F}$. Let μ be a finite measure on $\mathcal{E} \times \mathcal{F}$. Letting $\alpha(E) \equiv \mu(E \times Y)$, there exist probability measures ν_x unique up to a set of α measure zero such that for all $f \geq 0$ and $\mathcal{E} \times \mathcal{F}$ measurable,

$$\int_{X \times Y} f d\mu = \int_X \int_Y f d\nu_x d\alpha$$

Proof: Consider for $E, F \in \mathcal{E}, \mathcal{F}$ respectively, $\int_{X \times Y} \mathcal{X}_E(x) \mathcal{X}_F(y) d\mu$. Letting $\alpha(E) \equiv \mu(E \times Y)$, it follows that $E \rightarrow \int_{X \times Y} \mathcal{X}_E(x) \mathcal{X}_F(y) d\mu$ is absolutely continuous with respect to α . Therefore, there exists a unique nonnegative function in L^1 called h_F such that for any $E \in \mathcal{E}$,

$$\int_{X \times Y} \mathcal{X}_E(x) \mathcal{X}_F(y) d\mu = \int_X \mathcal{X}_E(x) h_F(x) d\alpha \quad (10.30)$$

That is h_F does not depend on $E \in \mathcal{E}$. Note also that 10.30 shows right away that $h_F(x) \leq 1$ a.e. Just let $F = Y$. Also, this shows that $h_Y(x) = 1$ for α a.e. x because from 10.30,

$$\int_{X \times Y} \mathcal{X}_E(x) d\mu = \mu(E \times Y) = \alpha(E) = \int_X \mathcal{X}_E(x) h_Y(x) d\alpha$$

where $h_Y(x) \leq 1$. Now let $E_m = [h_Y < 1 - \frac{1}{m}]$ so $\alpha(E_m) \leq (1 - \frac{1}{m}) \alpha(E)$. Then the above shows $\alpha(E_m) = 0$ and so, taking a union for $m \in \mathbb{N}$, yields that the set where h_Y is less than 1 has α measure zero.

Now $F \rightarrow \int_X \mathcal{X}_E(x) h_F(x) d\alpha$ is clearly a measure because $\int_{X \times Y} \mathcal{X}_E(x) \mathcal{X}_F(y) d\mu = \int_X \mathcal{X}_E(x) h_F(x) d\alpha$ implies that if $\{F_i\}$ are disjoint, then $h_{\cup_i F_i} = \sum_i h_{F_i}$ this by the uniqueness in the Radon Nikodym theorem. That is, for fixed x , $F \rightarrow h_F(x)$ is a measure ν_x . Since $h_Y(x) = 1$ for α a.e. x and $0 \leq h_F(x) \leq 1$, $h_Y(x) = 1$ α a.e., we can let ν_x be a probability measure for α a.e. x . Summarizing,

$$\int_{X \times Y} \mathcal{X}_E(x) \mathcal{X}_F(y) d\mu = \int_X \mathcal{X}_E(x) \int_Y \mathcal{X}_F(y) d\nu_x d\alpha = \int_X \int_Y \mathcal{X}_{E \times F}(x, y) d\nu_x d\alpha$$

If $\hat{\nu}_x$ also works, then it must equal ν_x for α a.e. x .

Let the π system \mathcal{K} consist of $E \times F$ where $E \in \mathcal{E}$ and $F \in \mathcal{F}$. Let \mathcal{G} be those sets A of $\mathcal{E} \times \mathcal{F} \equiv \sigma(\mathcal{K})$ such that $\int_{X \times Y} \mathcal{X}_A d\mu = \int_X \int_Y \mathcal{X}_A(x, y) d\nu_x d\alpha$. Then \mathcal{G} contains \mathcal{K} and is closed with respect to countable disjoint unions and complements, the latter coming from the observation that $X \times Y \in \mathcal{K}$ which allows the same kind of argument used in the above treatment of product measures. Therefore, by Dynkin's lemma, $\mathcal{G} = \sigma(\mathcal{K})$ and so, using approximation with simple functions and the monotone convergence theorem, we obtain that for any f which is $\mathcal{E} \times \mathcal{F} \equiv \sigma(\mathcal{K})$ measurable and nonnegative the iterated integrals make sense and $\int_{X \times Y} f d\mu = \int_X \int_Y f d\nu_x d\alpha$ ■

These measures ν_x are called slicing measures. They can be used to define what is meant by independent random variables in probability. This also shows that a given μ is a product measure exactly when the ν_x don't depend on x .

Consider now many spaces $\prod_{i=1}^n X_i$ where μ is a measure on $\mathcal{E} \equiv \prod_{i=1}^n \mathcal{E}_i$ where this denotes the product measurable sets from the (X_i, \mathcal{E}_i) . Then for f nonnegative and \mathcal{E} measurable,

$$\int_{X_1 \times \cdots \times X_n} f d\mu = \int_{X_1 \times \cdots \times X_{n-1}} \int_{X_n} f d\nu_{(x_1, \dots, x_{n-1})}(x_n) d\nu_{(x_1, \dots, x_{n-1})}$$

Here for $E \in \prod_{i=1}^{n-1} \mathcal{E}_i$, $\nu(E) \equiv \nu(E \times X_n)$. This ν is denoted as $\nu(x_1, \dots, x_{n-1})$. Then this equals

$$\begin{aligned} & \int_{X_1 \times \cdots \times X_{n-2}} \int_{X_{n-1}} \int_{X_n} f d\nu_{(x_1, \dots, x_{n-1})}(x_n) d\nu_{(x_1, \dots, x_{n-2})}(x_{n-1}) d\nu_{(x_1, \dots, x_{n-2})} \\ & \quad \vdots \\ & \int_{X_1} \cdots \int_{X_n} f d\nu_{(x_1, \dots, x_{n-1})}(x_n) d\nu_{(x_1, \dots, x_{n-2})}(x_{n-1}) \cdots d\nu_{x_1}(x_2) d\nu_{x_1}(x_1) \end{aligned}$$

where for $E \in \mathcal{E}_1$

Corollary 10.14.13 For $\mathcal{E} \equiv \prod_{i=1}^n \mathcal{E}_i$, and $f : \prod_{i=1}^n X_i \rightarrow [0, \infty]$ for (X_i, \mathcal{E}_i) a measurable space and for μ a finite probability measure on \mathcal{E} , meaning $\mu(\prod_{i=1}^n X_i) = 1$, there are probability measures as denoted below by ν with various subscripts such that

$$\int_{X_1 \times \cdots \times X_n} f d\mu = \int_{X_1} \cdots \int_{X_n} f d\nu_{(x_1, \dots, x_{n-1})}(x_n) d\nu_{(x_1, \dots, x_{n-2})}(x_{n-1}) \cdots d\nu_{x_1}(x_2) d\nu_{x_1}(x_1)$$

10.15 Jensen's Inequality

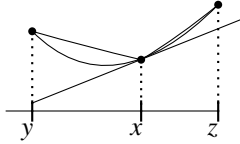
When you have $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is convex, then secant lines lie above the graph of ϕ . Say $x < w < z$ so $w = \lambda z + (1 - \lambda)x$ for some $\lambda \in (0, 1)$. Then referring to the following picture,

$$\frac{\phi(w) - \phi(x)}{w - x} \leq \frac{\lambda \phi(z) + (1 - \lambda) \phi(x) - \phi(x)}{(\lambda z + (1 - \lambda)x) - x} = \frac{\lambda (\phi(z) - \phi(x))}{\lambda (z - x)} = \frac{\phi(z) - \phi(x)}{z - x}$$

For $y < w < x$ so $w = \lambda y + (1 - \lambda)x$. Since $w - x < 0$,

$$\frac{\phi(w) - \phi(x)}{w - x} \geq \frac{\lambda \phi(y) + (1 - \lambda) \phi(x) - \phi(x)}{\lambda (y - x)} = \frac{\phi(y) - \phi(x)}{y - x}$$

Since x is arbitrary, this has shown that slopes of secant lines of the graph of ϕ over intervals increase as the intervals move to the right.



Lemma 10.15.1 *If $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is convex, then ϕ is continuous. Also, if ϕ is convex, $\mu(\Omega) = 1$, and $f, \phi(f) : \Omega \rightarrow \mathbb{R}$ are in $L^1(\Omega)$, then $\phi(\int_{\Omega} f d\mu) \leq \int_{\Omega} \phi(f) d\mu$.*

Proof: Let $\lambda \equiv \lim_{w \rightarrow x+} \frac{\phi(w) - \phi(x)}{w - x}$. Those slopes of secant lines are decreasing and so this limit exists. Then in the picture, for $w \in (x, z)$, $\phi(x) + \lambda(w - x) \leq \phi(w) \leq \phi(x) + \left(\frac{\phi(z) - \phi(x)}{z - x}\right)(w - x)$ and so ϕ is continuous from the right. A similar argument shows ϕ is continuous from the left. In particular, letting $\mu \equiv \lim_{w \rightarrow x-} \frac{\phi(x) - \phi(w)}{x - w} \leq \lambda$ because each of these slopes is smaller than the slopes whose inf gives λ . Then this shows that for $w \in (y, x)$, $\frac{\phi(w) - \phi(x)}{w - x} \leq \lambda$ so $\phi(w) - \phi(x) \geq \lambda(w - x)$ and so $\phi(w) \geq \phi(x) + \lambda(w - x)$ and for these w , $\frac{\phi(x) - \phi(w)}{x - w} \geq \frac{\phi(x) - \phi(y)}{x - y}$ so $\phi(w) \leq \phi(x) + \left(\frac{\phi(x) - \phi(y)}{x - y}\right)(w - x)$ so one obtains continuity from the left. This has also shown that for w not equal to x , $\phi(w) \geq \phi(x) + \lambda(w - x)$ or in other words, $\phi(x) \leq \phi(w) + \lambda(x - w)$. Letting $x = \int_{\Omega} f d\mu$, and using the λ whose existence was just established, for each ω ,

$$\phi\left(\int_{\Omega} f d\mu\right) \leq \phi(f(\omega)) + \lambda\left(\int_{\Omega} f d\mu - f(\omega)\right)$$

Do $\int_{\Omega} d\mu$ to both sides and use $\mu(\Omega) = 1$. Thus

$$\phi\left(\int_{\Omega} f d\mu\right) \leq \int_{\Omega} \phi(f) d\mu + \lambda\left(\int_{\Omega} f d\mu - \int_{\Omega} f d\mu\right) = \int_{\Omega} \phi(f) d\mu.$$

There are no difficulties with measurability because ϕ is continuous. ■

Corollary 10.15.2 *In the situation of Lemma 10.15.1 where $\mu(\Omega) = 1$, suppose f has values in $[0, \infty)$ and is measurable. Also suppose ϕ is convex and increasing on $[0, \infty)$. Then $\phi(\int_{\Omega} f d\mu) \leq \int_{\Omega} \phi(f) d\mu$.*

Proof: Let $f_n(\omega) = f(\omega)$ if $f(\omega) \leq n$ and let $f_n(\omega) = n$ if $f(\omega) \geq n$. Then both $f_n, \phi(f_n)$ are in $L^1(\Omega)$. Therefore, the above holds and $\phi(\int_{\Omega} f_n d\mu) \leq \int_{\Omega} \phi(f_n) d\mu$. Let $n \rightarrow \infty$ and use the monotone convergence theorem. ■

10.16 Faddeyev's Lemma

This next lemma is due to Faddeyev. I found it in [42].

Lemma 10.16.1 *Let f, g be nonnegative measurable nonnegative functions on a measure space (Ω, μ) . Then $\int f g d\mu = \int_0^\infty \int_{[g>t]} f d\mu dt = \int_0^\infty \int_0^\infty \mu([f>s] \cap [g>t]) ds dt$.*

Proof: First suppose $g = a\mathcal{X}_E$ where E is measurable, $a > 0$. Now $[g > t] = \emptyset$ if $t \geq a$ and it equals \mathcal{X}_E if $t < a$. Thus the right side equals $\int_0^a \int_E f d\mu dt = \int_0^a \int \mathcal{X}_E f d\mu = \int a \mathcal{X}_E f d\mu$ which equals the left side. Thus the first equation is true if $g = a\mathcal{X}_E$. Similar reasoning shows that when you have g a nonnegative simple function, $g = \sum_{i=1}^n a_i \mathcal{X}_{E_i}$ where we can arrange to have $\{a_i\}$ increasing, the first equation still holds. Now by the monotone convergence theorem, this yields the desired result for the first equation.

To get the second equal sign, note that

$$\begin{aligned} \int_0^\infty \int_{[g>t]} f d\mu dt &= \int_0^\infty \int \mathcal{X}_{[g>t]} f d\mu dt = \int_0^\infty \int_0^\infty \mu([\mathcal{X}_{[g>t]} f > s]) ds dt \\ &= \int_0^\infty \int_0^\infty \mu([f > s] \cap [g > t]) ds dt \blacksquare \end{aligned}$$

10.17 Exercises

1. Let $\Omega = \mathbb{N} = \{1, 2, \dots\}$. Let $\mathcal{F} = \mathcal{P}(\mathbb{N})$, the set of all subsets of \mathbb{N} , and let $\mu(S)$ = number of elements in S . Thus $\mu(\{1\}) = 1 = \mu(\{2\})$, $\mu(\{1, 2\}) = 2$, etc. In this case, all functions are measurable. For a nonnegative function, f defined on \mathbb{N} , show $\int_{\mathbb{N}} f d\mu = \sum_{k=1}^\infty f(k)$. What do the monotone convergence and dominated convergence theorems say about this example?
2. For the measure space of Problem 1, give an example of a sequence of nonnegative measurable functions $\{f_n\}$ converging pointwise to a function f , such that inequality is obtained in Fatou's lemma.
3. If $(\Omega, \mathcal{F}, \mu)$ is a measure space and $f \geq 0$ is measurable, show that if $g(\omega) = f(\omega)$ a.e. ω and $g \geq 0$, then $\int g d\mu = \int f d\mu$. Show that if $f, g \in L^1(\Omega)$ and $g(\omega) = f(\omega)$ a.e. then $\int g d\mu = \int f d\mu$.
4. Let $\{f_n\}, f$ be measurable functions with values in \mathbb{C} . $\{f_n\}$ converges in measure if $\lim_{n \rightarrow \infty} \mu(x \in \Omega : |f(x) - f_n(x)| \geq \varepsilon) = 0$ for each fixed $\varepsilon > 0$. Prove the theorem of F. Riesz. If f_n converges to f in measure, then there exists a subsequence $\{f_{n_k}\}$ which converges to f a.e. In case μ is a probability measure, this is called convergence in probability. It does not imply pointwise convergence but does imply that there is a subsequence which converges pointwise off a set of measure zero. **Hint:** Choose n_1 such that $\mu(x : |f(x) - f_{n_1}(x)| \geq 1) < 1/2$. Choose $n_2 > n_1$ such that $\mu(x : |f(x) - f_{n_2}(x)| \geq 1/2) < 1/2^2$. Choose $n_3 > n_2$ such that $\mu(x : |f(x) - f_{n_3}(x)| \geq 1/3) < 1/2^3$, etc. Now consider what it means for $f_{n_k}(x)$ to fail to converge to $f(x)$. Use the Borel Cantelli Lemma 9.2.5 on Page 243.
5. Let (X, \mathcal{F}, μ) be a regular measure space. For example, it could be \mathbb{R}^p with Lebesgue measure. Why do we care about a measure space being regular? This problem will show why. Suppose that closures of balls are compact as in the case of \mathbb{R}^p .

- (a) Let $\mu(E) < \infty$. By regularity, there exists $K \subseteq E \subseteq V$ where K is compact and V is open such that $\mu(V \setminus K) < \varepsilon$. Show there exists W open such that $K \subseteq \bar{W} \subseteq V$ and \bar{W} is compact. Now show there exists a function h such that h has values in $[0, 1]$, $h(x) = 1$ for $x \in K$, and $h(x)$ equals 0 off W . **Hint:** You might consider Problem 11 on Page 259.
- (b) Show that $\int |\mathcal{K}_E - h| d\mu < \varepsilon$
- (c) Next suppose $s = \sum_{i=1}^n c_i \mathcal{K}_{E_i}$ is a nonnegative simple function where each $\mu(E_i) < \infty$. Show there exists a continuous nonnegative function h which equals zero off some compact set such that $\int |s - h| d\mu < \varepsilon$
- (d) Now suppose $f \geq 0$ and $f \in L^1(\Omega)$. Show that there exists $h \geq 0$ which is continuous and equals zero off a compact set such that $\int |f - h| d\mu < \varepsilon$
- (e) If $f \in L^1(\Omega)$ with complex values, show the conclusion in the above part of this problem is the same.
6. Let $(\Omega, \mathcal{F}, \mu)$ be a measure space and suppose $f, g : \Omega \rightarrow (-\infty, \infty]$ are measurable. Prove the sets $\{\omega : f(\omega) < g(\omega)\}$ and $\{\omega : f(\omega) = g(\omega)\}$ are measurable. **Hint:** The easy way to do this is to write

$$\{\omega : f(\omega) < g(\omega)\} = \bigcup_{r \in \mathbb{Q}} [f < r] \cap [g > r].$$

Note that $l(x, y) = x - y$ is not continuous on $(-\infty, \infty]$ so the obvious idea doesn't work. Here $[g > r]$ signifies $\{\omega : g(\omega) > r\}$.

7. Let $\{f_n\}$ be a sequence of real or complex valued measurable functions. Let

$$S = \{\omega : \{f_n(\omega)\} \text{ converges}\}.$$

Show S is measurable. **Hint:** You might try to exhibit the set where f_n converges in terms of countable unions and intersections using the definition of a Cauchy sequence.

8. Suppose $u_n(t)$ is a differentiable function for $t \in (a, b)$ and suppose that for $t \in (a, b)$, $|u_n(t)|, |u'_n(t)| < K_n$ where $\sum_{n=1}^{\infty} K_n < \infty$. Show $(\sum_{n=1}^{\infty} u_n(t))' = \sum_{n=1}^{\infty} u'_n(t)$. **Hint:** This is an exercise in the use of the dominated convergence theorem and the mean value theorem.
9. Suppose $\{f_n\}$ is a sequence of nonnegative measurable functions defined on a measure space, $(\Omega, \mathcal{S}, \mu)$. Show that $\int \sum_{k=1}^{\infty} f_k d\mu = \sum_{k=1}^{\infty} \int f_k d\mu$. **Hint:** Use the monotone convergence theorem along with the fact the integral is linear.

10. Explain why for each $t > 0$, $x \rightarrow e^{-tx}$ is a function in $L^1(\mathbb{R})$ and $\int_0^{\infty} e^{-tx} dx = \frac{1}{t}$. Thus

$$\int_0^R \frac{\sin(t)}{t} dt = \int_0^R \int_0^{\infty} \sin(t) e^{-tx} dx dt$$

Now explain why you can change the order of integration in the above iterated integral. Then compute what you get. Next pass to a limit as $R \rightarrow \infty$ and show $\int_0^{\infty} \frac{\sin(t)}{t} dt = \frac{1}{2}\pi$. This is a very important integral. Note that the thing on the left is an improper integral. $\sin(t)/t$ is not Lebesgue integrable because it is not absolutely integrable. That is $\int_0^{\infty} \left| \frac{\sin t}{t} \right| dm = \infty$. It is important to understand that the Lebesgue theory of integration only applies to nonnegative functions and those which are absolutely integrable.

11. Let the rational numbers in $[0, 1]$ be $\{r_k\}_{k=1}^\infty$ and define

$$f_n(t) = \begin{cases} 1 & \text{if } t \in \{r_1, \dots, r_n\} \\ 0 & \text{if } t \notin \{r_1, \dots, r_n\} \end{cases}$$

Show that $\lim_{n \rightarrow \infty} f_n(t) = f(t)$ where f is one on the rational numbers and 0 on the irrational numbers. Explain why each f_n is Riemann integrable but f is not. However, each f_n is actually a simple function and its Lebesgue and Riemann integral is equal to 0. Apply the monotone convergence theorem to conclude that f is Lebesgue integrable and in fact, $\int f dm = 0$.

12. Show $\lim_{n \rightarrow \infty} \frac{n}{2^n} \sum_{k=1}^n \frac{2^k}{k} = 2$. This problem was shown to me by Shane Tang, a former student. It is a nice exercise in dominated convergence theorem if you massage it a little. **Hint:**

$$\frac{n}{2^n} \sum_{k=1}^n \frac{2^k}{k} = \sum_{k=1}^n 2^{k-n} \frac{n}{k} = \sum_{l=0}^{n-1} 2^{-l} \frac{n}{n-l} = \sum_{l=0}^{n-1} 2^{-l} \left(1 + \frac{l}{n-l}\right) \leq \sum_{l=0}^{n-1} 2^{-l} (1+l)$$

13. Suppose you have a real vector space $E(\omega)$ which is a subspace of a normed linear space V , this for each $\omega \in \Omega$ where (Ω, \mathcal{F}) is a measurable space. Suppose $E(\omega) = \text{span}(b_1(\omega), \dots, b_n(\omega))$ where these $b_i(\omega)$ are linearly independent and each is measurable into V . Define $\theta(\omega) : \mathbb{R}^n \rightarrow E(\omega)$ by $\theta(\omega)(\sum_{i=1}^n a_i e_i) \equiv \sum_{i=1}^n a_i b_i(\omega)$. Show that $\theta(\omega)$ maps functions measurable into \mathbb{R}^n to functions measurable into V . Now show $\theta(\omega)^{-1}$ also maps functions measurable into V to functions measurable into \mathbb{R}^n . **Hint:** For the second part you need to start with a function $\omega \rightarrow h(\omega)$ which is measurable into V with values in $E(\omega)$. Thus $h(\omega) = \sum_{i=1}^n a_i(\omega) b_i(\omega)$. You need to verify that the a_i are measurable. To do this, assume first that $\|h(\omega)\|$ is bounded by some constant M . Then consider $S_r \equiv \{\omega : \inf_{|a| > r} \|\sum_i a_i b_i(\omega)\| > M\}$. Explain why every ω is in some S_r . Then consider $\Phi(a, \omega)$ which will be defined as $-\|\sum_i a_i b_i(\omega) - h(\omega)\|$ for $\omega \in S_r$. Thus the maximum of this functional for $\omega \in \Omega$ is 0. Show that for $\omega \in S_r$ the maximum will occur on the set $|a| \leq M+1$. Then apply Kuratowski's lemma. Finally consider a truncation of h called h_m and apply what was just shown to this truncation which has $\|h_m(\omega)\| \leq m$. Then let $m \rightarrow \infty$ and observe that for large enough m , $h_m(\omega) = h(\omega)$ and so the $a_i^m(\omega)$ are also constant from this value onward. Thus $\lim_{m \rightarrow \infty} a_i^m(\omega) \equiv a_i(\omega)$ exists and a_i is measurable, being the limit of measurable functions.
14. Give an example of a sequence of functions $\{f_n\}$, $f_n \geq 0$ and a function $f \geq 0$ such that $f(x) = \liminf_{n \rightarrow \infty} f_n(x)$ but $\int f dm < \liminf_{n \rightarrow \infty} \int f_n dm$ so you get strict inequality in Fatou's lemma.
15. Let f be a nonnegative Riemann integrable function defined on $[a, b]$. Thus there is a unique number between all the upper sums and lower sums. First explain why, if $a_i \geq 0$, $\int \sum_{i=1}^n a_i \mathcal{X}_{[t_i, t_{i+1})}(t) dm = \sum_i a_i (t_i - t_{i-1})$. Explain why there exists an increasing sequence of Borel measurable functions $\{g_n\}$ converging to a Borel measurable function g , and a decreasing sequence of functions $\{h_n\}$ which are also Borel measurable converging to a Borel measurable function h such that $g_n \leq f \leq h_n$,

$$\int g_n dm \text{ equals a lower sum, } \int h_n dm \text{ equals an upper sum}$$

and $\int (h - g) dm = 0$. Explain why $\{x : f(x) \neq g(x)\}$ is a set of measure zero. Then explain why f is measurable and $\int_a^b f(x) dx = \int f dm$ so that the Riemann integral gives the same answer as the Lebesgue integral. Here m is one dimensional Lebesgue measure discussed earlier.

16. Let λ, μ be finite measures. We say $\lambda \ll \mu$ if whenever $\mu(E) = 0$ it follows $\lambda(E) = 0$. Show that if $\lambda \ll \mu$, then for every $\varepsilon > 0$ there exists $\delta > 0$ such that if $\mu(E) < \delta$, then $\lambda(E) < \varepsilon$.
17. If λ is a signed measure with values in \mathbb{R} so that when $\{E_i\}$ are disjoint, $\sum_i \lambda(E_i)$ converges, show that the infinite series converges absolutely also.
18. In the Radon Nikodym Theorem 10.13.7, show that if f, \hat{f} both work, then $f = \hat{f}$ a.e.
19. Suppose $\nu \ll \mu$ where these are finite measures so there exists $h \geq 0$ and measurable such that $\nu(E) = \int_E h d\mu$ by the Radon Nikodym theorem. Show that if f is measurable and non-negative, then $\int f d\nu = \int f h d\mu$. **Hint:** It holds if f is χ_E and so it holds for a simple function. Now consider a sequence of simple functions increasing to f and use the monotone convergence theorem.
20. If the functions f_i of the above problem are “independent” you have

$$\mu(\cap_{i=1}^m [f_i \geq s_i]) = \prod_{i=1}^m \mu([f_i \geq s_i])$$

Suppose then that $\{f_i\}_{i=1}^m$ are independent. Show $\int_{\Omega} \prod_{i=1}^m f_i d\mu = \prod_{i=1}^m \int_{\Omega} f_i d\mu$. If μ is a probability measure, then such measurable functions are called random variables.

21. Suppose the situation of Corollary 10.14.13 in which μ is a probability measure on $(\prod_{i=1}^n X_i, \mathcal{E})$ where \mathcal{E} consists of the product measurable sets and for f a nonnegative \mathcal{E} measurable function,

$$\begin{aligned} & \int_{X_1 \times \cdots \times X_n} f d\mu \\ &= \int_{X_1} \cdots \int_{X_n} f d\nu_{(x_1, \dots, x_{n-1})}(x_n) d\nu_{(x_1, \dots, x_{n-2})}(x_{n-1}) \cdots d\nu_{x_1}(x_2) d\nu(x_1) \end{aligned}$$

Show that if the slicing measures do not depend on the subscripts, then whenever $E_{k+1} \in \mathcal{E}_{k+1}$,

$$\nu_{(x_1, \dots, x_k)}(E_{k+1}) = \mu(X_1 \times \cdots \times X_k \times E_{k+1} \times X_{k+2} \times \cdots \times X_n) \equiv \mu(E_{k+1})$$

where $E_{k+1} \rightarrow \mu(E_{k+1})$ is a probability measure which does not depend on the vector (x_1, \dots, x_k) . If any of the slicing measures does depend on the subscripts, show that something like this cannot take place.

Hint: Consider $\mathcal{X}_{E_{k+1}} = f$.

22. Suppose $\nu_{(x_1, \dots, x_k)} = \nu_k$ aside from an appropriate set of $\nu_{(x_1, \dots, x_k)}$ measure zero, where ν_k does not depend on (x_1, \dots, x_k) . Show that then, $\int_{X_1 \times \cdots \times X_n} \prod_{i=1}^n \mathcal{X}_{E_i} d\mu = \prod_{i=1}^n \nu_i(E_i) = \prod_{i=1}^n \mu(E_i)$. This has to do with the notion of independent events.

23. Use Jensen's inequality in Corollary 10.15.2 to show that if f is nonnegative and measurable, then for $p > 1$ show that whenever μ is a finite measure, then if $f^p \in L^1(\Omega)$ it follows that $f \in L^1(\Omega)$. Give an example to show that this is not necessarily true if $\mu(\Omega) = \infty$. **Hint:** For the second part, you might consider $\Omega = \mathbb{N}$, the σ algebra the set of all subsets, and $\mu(S)$ equal to the number of elements in S . Maybe $f(n) = 1/n$.

Chapter 11

Regular Measures

So far, most examples have been in one dimensional settings. This is about to change.

11.1 Regular Measures in a Metric Space

In this section X will be a metric space in which the closed balls are compact. The extra generality involving a metric space instead of \mathbb{R}^p would allow the consideration of manifolds for example. However, \mathbb{R}^p is an important case.

Definition 11.1.1 *The symbol $C_c(V)$ for V an open set will denote the continuous functions having compact support which is contained in V . Recall that the support of a continuous function f is defined as the closure of the set on which the function is nonzero. $L: C_c(X) \rightarrow \mathbb{C}$ is called a positive linear functional if it is linear, $L(\alpha f + \beta g) = \alpha Lf + \beta Lg$ and satisfies $Lf \leq Lg$ if $f \leq g$. Also, recall that a measure μ is regular on some σ algebra \mathcal{F} containing the Borel sets if for every $F \in \mathcal{F}$,*

$$\begin{aligned}\mu(F) &= \sup \{ \mu(K) : K \subseteq F \text{ and } K \text{ compact} \} \\ \mu(F) &= \inf \{ \mu(V) : V \supseteq F \text{ and } V \text{ is open} \}\end{aligned}$$

A complete measure, finite on compact sets, which is regular as above, is called a Radon measure. A set is called an F_σ set if it is the countable union of closed sets and a set is G_δ if it is the countable intersection of open sets.

Remarkable things happen in the above context. Some are described in the following proposition.

Proposition 11.1.2 *Suppose (X, d) is a metric space in which the closed balls are compact and X is a countable union of closed balls. Also suppose (X, \mathcal{F}, μ) is a complete measure space, \mathcal{F} contains the Borel sets, and that μ is regular and finite on measurable subsets of finite balls. Then*

1. *For each $E \in \mathcal{F}$, there is an F_σ set F and a G_δ set G such that $F \subseteq E \subseteq G$ and $\mu(G \setminus F) = 0$.*
2. *Also if $f \geq 0$ is \mathcal{F} measurable, then there exists $g \leq f$ such that g is Borel measurable and $g = f$ a.e. and $h \geq f$ such that h is Borel measurable and $h = f$ a.e.*
3. *If $E \in \mathcal{F}$ is a bounded set contained in a ball $B(x_0, r) = V$, then there exists a sequence of continuous functions in $C_c(V)$ $\{h_n\}$ having values in $[0, 1]$ and a set of measure zero N such that for $x \notin N$, $h_n(x) \rightarrow \chi_E(x)$. Also $\int |h_n - \chi_E| d\mu \rightarrow 0$. Letting \tilde{N} be a G_δ set of measure zero containing N , $h_n \chi_{\tilde{N}^c} \rightarrow \chi_F$ where $F \subseteq E$ and $\mu(E \setminus F) = 0$.*
4. *If $f \in L^1(X, \mathcal{F}, \mu)$, there exists $g \in C_c(X)$, such that $\int_X |f - g| d\mu < \epsilon$. There also exists a sequence of functions in $C_c(X)$ $\{g_n\}$ which converges pointwise to f .*

Proof: 1. Let $R_n \equiv B(x_0, n)$, $R_0 = \emptyset$. If E is measurable, let $E_n \equiv E \cap (R_n \setminus R_{n-1})$. Thus these E_n are disjoint and their union is E . By outer regularity, there exists open $U_n \supseteq E_n$ such that $\mu(U_n \setminus E_n) < \varepsilon/2^n$. Now if $U \equiv \bigcup_n U_n$, it follows that $U \setminus E \subseteq \bigcup_n (U_n \setminus E_n)$ so $\mu(U \setminus E) \leq \sum_{n=1}^{\infty} \frac{\varepsilon}{2^n} = \varepsilon$. This has shown that there exists an open set U containing E such that $\mu(U \setminus E) \leq \varepsilon$. Let V_n be open, containing E and $\mu(V_n \setminus E) < \frac{1}{2^n}$, $V_n \supseteq V_{n+1}$. Let $G \equiv \bigcap_n V_n$. This is a G_δ set containing E and $\mu(G \setminus E) \leq \mu(V_n \setminus E) < \frac{1}{2^n}$ and so $\mu(G \setminus E) = 0$. By inner regularity, there is F_n an F_σ set contained in E_n with $\mu(E_n \setminus F_n) = 0$. Then let $F \equiv \bigcup_n F_n$. This F is an F_σ set and $\mu(E \setminus F) \leq \sum_n \mu(E_n \setminus F_n) = 0$. Thus $F \subseteq E \subseteq G$ and $\mu(G \setminus F) \leq \mu(G \setminus E) + \mu(E \setminus F) = 0$.

2. If f is measurable and nonnegative, from Theorem 9.1.6 there is an increasing sequence of simple functions s_n such that $\lim_{n \rightarrow \infty} s_n(x) = f(x)$. Say $s_n(x) \equiv \sum_{k=1}^{m_n} c_k^n \mathcal{X}_{E_k^n}(x)$. Let $m_p(E_k^n \setminus F_k^n) = 0$ where F_k^n is an F_σ set. Replace E_k^n with F_k^n and let \tilde{s}_n be the resulting simple function. Let $g(x) \equiv \lim_{n \rightarrow \infty} \tilde{s}_n(x)$. Then g is Borel measurable and $g \leq f$ and $g = f$ except for a set of measure zero, the union of the sets where s_n is not equal to \tilde{s}_n . As to the other claim, let $h_n(x) \equiv \sum_{k=1}^{\infty} \mathcal{X}_{A_{kn}}(x) \frac{k}{2^n}$ where A_{kn} is a G_δ set containing $f^{-1}((\frac{k-1}{2^n}, \frac{k}{2^n}])$ for which $\mu(A_{kn} \setminus f^{-1}((\frac{k-1}{2^n}, \frac{k}{2^n}])) \equiv \mu(D_{kn}) = 0$. If $N = \bigcup_{k,n} D_{kn}$, then N is a set of measure zero. On N^C , $h_n(x) \rightarrow f(x)$. Let $h(x) = \liminf_{n \rightarrow \infty} h_n(x)$. Note that $\mathcal{X}_{A_{kn}}(x) \frac{k}{2^n} \geq \mathcal{X}_{f^{-1}((\frac{k-1}{2^n}, \frac{k}{2^n}])}(x) \frac{k}{2^n}$ and so $h_n(x) \geq h(x)$ and $\liminf_{n \rightarrow \infty} h_n(x)$ is Borel measurable because each h_n is.

3. Let $K_n \subseteq E \subseteq V_n$ with K_n compact and V_n open such that $V_n \subseteq B(x_0, r)$ and also that $\mu(V_n \setminus K_n) < 2^{-(n+1)}$. Then from Lemma 3.12.4, there is h_n with $K_n \prec h_n \prec V_n$. Then $\int |h_n - \mathcal{X}_E| d\mu < 2^{-n}$ and so

$$\mu\left(|h_n - \mathcal{X}_E| > \left(\frac{2}{3}\right)^n\right) < \left(\left(\frac{3}{2}\right)^n \int_{[|h_n - \mathcal{X}_E| > (\frac{2}{3})^n]} |h_n - \mathcal{X}_E| d\mu\right) \leq \left(\frac{3}{4}\right)^n$$

By Lemma 9.2.5 there is a set of measure zero N such that if $x \notin N$, it is in only finitely many of the sets $[|h_n - \mathcal{X}_E| > (\frac{2}{3})^n]$. Thus on N^C , eventually, for all k large enough, $|h_k - \mathcal{X}_E| \leq (\frac{2}{3})^k$ so $h_k(x) \rightarrow \mathcal{X}_E(x)$ off N . The assertion about convergence of the integrals follows from the dominated convergence theorem and the fact that each h_n is nonnegative, bounded by 1, ($K_n \prec h_n \prec V_n$) and is 0 off some ball. In the last claim, it only remains to verify that $h_n \mathcal{X}_{\tilde{N}^C}$ converges to an indicator function because each $h_n \mathcal{X}_{\tilde{N}^C}$ is Borel measurable. ($\tilde{N} \supseteq N$ and \tilde{N} is a Borel set and $\mu(\tilde{N} \setminus N) = 0$) Thus its limit will also be Borel measurable. However, $h_n \mathcal{X}_{\tilde{N}^C}$ converges to 1 on $E \cap \tilde{N}^C$, 0 on $E^C \cap \tilde{N}^C$ and 0 on \tilde{N} . Thus $E \cap \tilde{N}^C = F$ and $h_n \mathcal{X}_{\tilde{N}^C}(x) \rightarrow \mathcal{X}_F$ where $F \subseteq E$ and $\mu(E \setminus F) \leq \mu(\tilde{N}) = 0$.

4. It suffices to assume $f \geq 0$ because you can consider the positive and negative parts of the real and imaginary parts of f and reduce to this case. Let $f_n(x) \equiv \mathcal{X}_{B(x_0, n)}(x) f(x)$. Then by the dominated convergence theorem, if n is large enough, $\int |f - f_n| d\mu < \varepsilon$. There is a nonnegative simple function $s \leq f_n$ such that $\int |f_n - s| d\mu < \varepsilon$. This follows from picking k large enough in an increasing sequence of simple functions $\{s_k\}$ converging to f_n and the dominated convergence theorem. Say $s(x) = \sum_{k=1}^m c_k \mathcal{X}_{E_k}(x)$. Then let $K_k \subseteq E_k \subseteq V_k$ where K_k, V_k are compact and open respectively and $\sum_{k=1}^m c_k \mu(V_k \setminus K_k) < \varepsilon$. By Lemma

3.12.4, there exists h_k with $K_k \prec h_k \prec V_k$. Then

$$\begin{aligned} \int \left| \sum_{k=1}^m c_k \mathcal{X}_{E_k}(x) - \sum_{k=1}^m c_k h_k(x) \right| d\mu &\leq \sum_k c_k \int |\mathcal{X}_{E_k}(x) - h_k(x)| dx \\ &< 2 \sum_k c_k \mu(V_k \setminus K_k) < 2\varepsilon \end{aligned}$$

Let $g \equiv \sum_{k=1}^m c_k h_k(x)$. Thus $\int |s - g| d\mu \leq 2\varepsilon$. Then

$$\int |f - g| d\mu \leq \int |f - f_n| d\mu + \int |f_n - s| d\mu + \int |s - g| d\mu < 4\varepsilon$$

Since ε is arbitrary, this proves the first part of 4. For the second part, let $g_n \in C_c(X)$ such that $\int |f - g_n| d\mu < 2^{-n}$. Let $A_n \equiv \left\{ x : |f - g_n| > \left(\frac{2}{3}\right)^n \right\}$. Then

$$\mu(A_n) \leq \left(\frac{3}{2}\right)^n \int_{A_n} |f - g_n| d\mu \leq \left(\frac{3}{4}\right)^n.$$

Thus, if N is all x in infinitely many A_n , then by the Borel Cantelli lemma, $\mu(N) = 0$ and if $x \notin N$, then x is in only finitely many A_n and so for all n large enough, $|f(x) - g_n(x)| \leq \left(\frac{2}{3}\right)^n$. ■

For the rest of this chapter, I will specialize to \mathbb{R}^p or at least a finite dimensional normed linear space. For different proofs and some results which are not discussed here, a good source is [17] which is where I first read some of these things.

Recall the following Besicovitch covering theorem for Radon measures. It is Corollary 9.12.3 on Page 264 and the earlier version, Theorem 4.5.8 on Page 119 which are listed here for the sake of convenience.

Corollary 11.1.3 *Let E be a nonempty set and let μ be a Radon measure on a σ algebra which contains the Borel sets of \mathbb{R}^p . Suppose \mathcal{F} is a collection of closed balls which cover E in the sense of Vitali. Then there exists a sequence of disjoint closed balls $\{B_i\} \subseteq \mathcal{F}$ such that $\bar{\mu}\left(E \setminus \bigcup_{j=1}^N B_j\right) = 0, N \leq \infty$.*

Theorem 11.1.4 *There exists a constant N_p , depending only on p with the following property. If \mathcal{F} is any collection of nonempty balls in \mathbb{R}^p with*

$$\sup \{\text{diam}(B) : B \in \mathcal{F}\} = D < \infty$$

and if A is the set of centers of the balls in \mathcal{F} , then there exist subsets of \mathcal{F} , $\mathcal{H}_1, \dots, \mathcal{H}_{N_p}$, such that each \mathcal{H}_i is a countable collection of disjoint balls from \mathcal{F} (possibly empty) and

$$A \subseteq \bigcup_{i=1}^{N_p} \mathcal{H}_i \cup \{B : B \in \mathcal{H}_i\}.$$

11.2 Constructing Measures from Functionals

Here is a theorem which is the main result on measures and functionals defined on a space of continuous functions. The typical situation is of a metric space in which closed balls are compact.

Definition 11.2.1 $C_c(X)$ will denote the complex valued functions which have compact support in some metric space X . This is clearly a linear space. Then a linear function $L : C_c(X) \rightarrow \mathbb{C}$ is called “positive” if whenever $f \geq 0$, then $Lf \geq 0$.

Theorem 11.2.2 Let $L : C_c(X) \rightarrow \mathbb{C}$ be a positive linear functional where X is a metric space and X is a countable union of compact sets. Then there exists a complete measure μ defined on a σ algebra \mathcal{F} which contains the Borel sets $\mathcal{B}(X)$ which is finite on compact sets and has the following properties.

1. μ is regular and if E is measurable, there are F_σ and G_δ sets F, G such that $F \subseteq E \subseteq G$ and $\mu(G \setminus F) = 0$.
2. For all $f \in C_c(X)$, $Lf = \int_X f d\mu$

Proof: See the notation and lemmas near Definition 3.12.3 having to do with partitions of unity on a metric space for what is needed in this proof. For V open, let $\bar{\mu}(V) \equiv \sup\{Lf : f \prec V\}$. Then for an arbitrary set F , let $\bar{\mu}(F) \equiv \inf\{\bar{\mu}(V) : V \supseteq F\}$, $\bar{\mu}(\emptyset) \equiv 0$. In what follows, V will be an open set and K a compact set.

Claim 1: $\bar{\mu}$ is well defined.

Proof of Claim 1: Note there are two descriptions of $\bar{\mu}(V)$ for V open. They need to be the same. Let $\bar{\mu}_1$ be the definition involving supremums of Lf and let $\bar{\mu}$ be the general definition. Let $V \subseteq U$ where V, U open. Then by definition, $\bar{\mu}(V) \leq \bar{\mu}_1(U)$ and so $\bar{\mu}(V) \equiv \inf\{\bar{\mu}_1(U) : U \supseteq V\} \geq \bar{\mu}_1(V)$. However, $V \subseteq V$ and so $\bar{\mu}(V) \leq \bar{\mu}_1(V)$. ■

Claim 2: $\bar{\mu}$ is finite on compact sets. Also, if $K \prec f$, it follows that $\bar{\mu}(K) \leq L(f) < \infty$.

Proof of Claim 2: Let $K \prec f \prec X$. Let $V_\varepsilon \equiv \{x : f(x) > 1 - \varepsilon\}$, an open set since f is continuous. Then let $g \prec V_\varepsilon$ so it follows that $\frac{f}{1-\varepsilon} \geq g$. Then $L(g) \leq \frac{1}{1-\varepsilon} L(f) < \infty$. Then taking the sup over all such g , it follows that $\bar{\mu}(K) \leq \bar{\mu}(V_\varepsilon) \leq \frac{1}{1-\varepsilon} Lf$. Now let $\varepsilon \rightarrow 0$ and conclude that $\bar{\mu}(K) \leq L(f)$. ■

Claim 3: $\bar{\mu}$ is subadditive: $\bar{\mu}(\cup_i E_i) \leq \sum_i \bar{\mu}(E_i)$.

Proof of Claim 3: First consider the case of open sets. Let $V = \cup_i V_i$. Let $l < \bar{\mu}(V)$. Then there exists $f \prec V$ with $Lf > l$. Then $\sup(f)$ is contained in $\cup_{i=1}^n V_i$. Now let $\sup \psi_i \subseteq V_i$ and $\sum_{i=1}^n \psi_i = 1$ on $\sup(f)$. This is from Theorem 3.12.5. Then

$$l < Lf = \sum_{i=1}^n L(\psi_i f) \leq \sum_{i=1}^n \bar{\mu}(V_i) \leq \sum_i \bar{\mu}(V_i).$$

Since l is arbitrary, it follows that $\bar{\mu}(V) \leq \sum_i \bar{\mu}(V_i)$. Now consider the general case. Let $E = \cup_i E_i$. If $\sum_i \bar{\mu}(E_i) = \infty$, there is nothing to show. Assume then that this sum is finite and let $V_i \supseteq E_i$, $\bar{\mu}(E_i) + \frac{\varepsilon}{2^i} > \bar{\mu}(V_i)$. Then

$$\bar{\mu}(E) \leq \bar{\mu}(\cup_i V_i) \leq \sum_i \bar{\mu}(V_i) \leq \sum_i \left(\bar{\mu}(E_i) + \frac{\varepsilon}{2^i} \right) = \sum_i \bar{\mu}(E_i) + \varepsilon$$

Since ε is arbitrary, this shows $\bar{\mu}$ is subadditive. ■

Claim 4: If $\text{dist}(A, B) = \delta > 0$, then $\bar{\mu}(A \cup B) = \bar{\mu}(A) + \bar{\mu}(B)$.

Proof of Claim 4: If the right side is infinite, there is nothing to show so we can assume that $\bar{\mu}(A), \bar{\mu}(B)$ are both finite. First suppose U, V are open and disjoint having finite outer measure. Let $\bar{\mu}(U) \leq Lf_1 + \varepsilon$ where $f_1 \prec U$ and let $f_2 \prec V$ with $\bar{\mu}(V) \leq L(f_2) + \varepsilon$. Then

$$\bar{\mu}(U \cup V) \leq \bar{\mu}(U) + \bar{\mu}(V) \leq Lf_1 + Lf_2 + 2\varepsilon \leq L(f_1 + f_2) + 2\varepsilon \leq \bar{\mu}(U \cup V) + 2\varepsilon$$

Since ε is arbitrary, this shows that $\bar{\mu}(U \cup V) = \bar{\mu}(U) + \bar{\mu}(V)$. Now in case A, B are as assumed, let $U \equiv \cup_{x \in U} B(x, \delta/3), V \equiv \cup_{x \in V} B(x, \delta/3)$. Then these are disjoint open sets containing A and B respectively. Then there is O open, $O \supseteq A \cup B$ such that $\bar{\mu}(A \cup B) + \varepsilon > \bar{\mu}(O)$. Replacing U with $U \subseteq O$ and V with $V \cap O$, we can assume $\bar{\mu}(A \cup B) + \varepsilon > \bar{\mu}(U \cup V)$. Then

$$\begin{aligned} \bar{\mu}(A) + \bar{\mu}(B) &\leq \bar{\mu}(U) + \bar{\mu}(V) = \bar{\mu}(U \cup V) \\ &< \varepsilon + \bar{\mu}(A \cup B) \leq \varepsilon + \bar{\mu}(A) + \bar{\mu}(B) \end{aligned}$$

Since ε is arbitrary, this shows that $\bar{\mu}(A) + \bar{\mu}(B) = \bar{\mu}(A \cup B)$.

From Theorem 9.5.4 there is a complete measure μ defined on a σ algebra \mathcal{F} which equals $\bar{\mu}$ on \mathcal{F} . From Claim 4 and Theorem 9.6.1, \mathcal{F} contains the Borel sets $\mathcal{B}(X)$. From the definition, μ is outer regular and so it follows from Theorem 9.8.6 that μ is regular because it is finite on compact sets and X is the union of countably many compact sets so μ is σ finite. Hence, by that theorem, the claimed approximation result also holds.

It only remains to show that the integrals with respect to the measure represent the functional. This will complete the proof.

Claim 5: $\int f d\mu = Lf$ for all $f \in C_c(X)$.

Proof: Let $f \in C_c(X)$, f real-valued, and suppose $f(X) \subseteq [a, b]$. Choose $t_0 < a$ and let $t_0 < t_1 < \dots < t_n = b$, $t_i - t_{i-1} < \varepsilon$. Let

$$E_i = f^{-1}((t_{i-1}, t_i]) \cap \text{spt}(f). \quad (11.1)$$

Note that $\cup_{i=1}^n E_i$ is a closed set equal to $\text{spt}(f)$.

$$\cup_{i=1}^n E_i = \text{spt}(f) \quad (11.2)$$

Since $X = \cup_{i=1}^n f^{-1}((t_{i-1}, t_i])$. Let $V_i \supseteq E_i$, V_i is open and let V_i satisfy

$$f(x) < t_i + \varepsilon \text{ for all } x \in V_i, \mu(V_i \setminus E_i) < \varepsilon/n. \quad (11.3)$$

By Theorem 3.12.5, there exists $h_i \in C_c(X)$ such that

$$h_i \prec V_i, \sum_{i=1}^n h_i(x) = 1 \text{ on } \text{spt}(f).$$

Now note that for each i , $f(x)h_i(x) \leq h_i(x)(t_i + \varepsilon)$. Therefore,

$$\begin{aligned} Lf &= L\left(\sum_{i=1}^n f h_i\right) \leq L\left(\sum_{i=1}^n h_i(t_i + \varepsilon)\right) = \sum_{i=1}^n (t_i + \varepsilon)L(h_i) \\ &= \sum_{i=1}^n (|t_0| + t_i + \varepsilon)L(h_i) - |t_0|L\left(\sum_{i=1}^n h_i\right). \end{aligned}$$

Now note that $|t_0| + t_i + \varepsilon \geq 0$ and so from the definition of μ and **claim 2**, this is no larger than

$$\sum_{i=1}^n (|t_0| + t_i + \varepsilon)\mu(V_i) - |t_0|\mu(\text{spt}(f)) \leq \sum_{i=1}^n (|t_0| + t_i + \varepsilon)(\mu(E_i) + \varepsilon/n) - |t_0|\mu(\text{spt}(f))$$

$$\begin{aligned}
& \leq |t_0| \overbrace{\sum_{i=1}^n \mu(E_i)}^{\mu(\text{spt}(f))} + \frac{\varepsilon}{n} n |t_0| + \sum_i t_i \mu(E_i) + \sum_i t_i \frac{\varepsilon}{n} + \sum_i \varepsilon \mu(E_i) + \frac{\varepsilon^2}{n} - |t_0| \mu(\text{spt}(f)) \\
& \leq \varepsilon |t_0| + \varepsilon (|t_0| + |b|) + \varepsilon \mu(\text{spt}(f)) + \varepsilon^2 + \sum_i t_i \mu(E_i) \\
& \leq \varepsilon |t_0| + \varepsilon (|t_0| + |b|) + 2\varepsilon \mu(\text{spt}(f)) + \varepsilon^2 + \sum_{i=1}^n t_{i-1} \mu(E_i) \\
& \leq \varepsilon (2|t_0| + |b| + 2\mu(\text{spt}(f)) + \varepsilon) + \int f d\mu
\end{aligned}$$

Since $\varepsilon > 0$ is arbitrary, $Lf \leq \int f d\mu$ for all $f \in C_c(X)$, f real. Hence equality holds because $L(-f) \leq -\int f d\mu$ so $L(f) \geq \int f d\mu$. Thus $Lf = \int f d\mu$ for all $f \in C_c(X)$. Just apply the result for real functions to the real and imaginary parts of f . ■

Using Corollary 9.8.9 we obtain the following corollary. Note that the conditions of the above theorem imply that X is a Polish space in the usual case where closed balls are compact.

Corollary 11.2.3 *If X is a Polish space then in the above theorem, we obtain inner regularity of μ in terms of compact sets. That is if $F \in \mathcal{F}$, then*

$$\mu(F) = \sup \{ \mu(K) : K \subseteq F, K \text{ compact} \}$$

11.3 The p Dimensional Lebesgue Measure

Theorem 11.2.2 will provide many examples of Radon measures on \mathbb{R}^p . Lebesgue measure is obtained by letting

$$Lf \equiv \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} f(x_1, \dots, x_p) dm_1(x_1) \cdots dm_p(x_p)$$

for $f \in C_c(\mathbb{R}^p)$. Thus Lebesgue measure is a Radon measure, denoted as m_p . In this case, the σ algebra will be denoted as \mathcal{F}_p . Lebesgue measure also has other very important properties. Integrals can be easily computed and the measure is translation invariant.

Theorem 11.3.1 *Whenever f is measurable and nonnegative, then whenever g is Borel measurable and equals f a.e. and h is Borel and equals f a.e.*

$$\begin{aligned}
& \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} h(x_1, \dots, x_p) dm_1(x_{i_1}) \cdots dm_p(x_{i_p}) = \\
& \int_{\mathbb{R}^p} f dm_p = \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} g(x_1, \dots, x_p) dm_1(x_{i_1}) \cdots dm_p(x_{i_p})
\end{aligned}$$

where (i_1, i_2, \dots, i_p) is any permutation of the integers $\{1, 2, \dots, p\}$. Also, m_p is regular and complete. If R is of the form $\prod_{i=1}^p I_i$ where I_i is an interval, then $m_p(R)$ is the product of the lengths of the sides of R . Also if $E \in \mathcal{F}_p$, then $m_p(x + E) = m_p(E)$.

Proof: Let \mathcal{H} consist of all open rectangles $\prod_i (a_i, b_i)$ along with \emptyset and \mathbb{R}^p . Thus this is a π system. Let $R_n \equiv \prod_{i=1}^p (-n, n)$. Let \mathcal{G} consist of the Borel sets $E \subseteq \mathbb{R}^p$ such that $\int_{R_n \cap E} dm_p = \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} \mathcal{X}_{R_n \cap E} dm_1(x_{i_1}) \cdots dm_p(x_{i_p})$. Then $\mathcal{H} \subseteq \mathcal{G}$. It is obvious from the monotone convergence theorem that \mathcal{G} is closed with respect to countable disjoint unions. Indeed, this theorem implies that for $E = \cup_i E_i$, the E_i disjoint,

$$\begin{aligned} \int_{R_n \cap E} dm_p &= \lim_{m \rightarrow \infty} \int_{R_n \cap \cup_{i=1}^m E_i} dm_p = \lim_{m \rightarrow \infty} \left(\sum_{i=1}^m \int \mathcal{X}_{R_n \cap E_i} dm_p \right) \\ &= \lim_{m \rightarrow \infty} \left(\sum_{i=1}^m \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} \mathcal{X}_{R_n \cap E_i} dm_1(x_{i_1}) \cdots dm_p(x_{i_p}) \right) \\ &= \lim_{m \rightarrow \infty} \left(\int_{\mathbb{R}} \cdots \int_{\mathbb{R}} \sum_{i=1}^m \mathcal{X}_{R_n \cap E_i} dm_1(x_{i_1}) \cdots dm_p(x_{i_p}) \right) \\ &= \lim_{m \rightarrow \infty} \left(\int_{\mathbb{R}} \cdots \int_{\mathbb{R}} \mathcal{X}_{R_n \cap \cup_{i=1}^m E_i} dm_1(x_{i_1}) \cdots dm_p(x_{i_p}) \right) \\ &= \left(\int_{\mathbb{R}} \cdots \int_{\mathbb{R}} \mathcal{X}_{R_n \cap E} dm_1(x_{i_1}) \cdots dm_p(x_{i_p}) \right) \end{aligned}$$

As to complements, $\int_{R_n} dm_p = \int_{R_n \cap E^c} dm_p + \int_{R_n \cap E} dm_p$. Thus

$$\begin{aligned} \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} \mathcal{X}_{R_n \cap E^c} dm_1(x_{i_1}) \cdots dm_p(x_{i_p}) &= \\ \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} (\mathcal{X}_{R_n} - \mathcal{X}_{R_n \cap E}) dm_1(x_{i_1}) \cdots dm_p(x_{i_p}) &= \int_{R_n \cap E^c} dm_p \end{aligned}$$

It follows that $\mathcal{G} = \mathcal{B}(\mathbb{R}^p)$, the Borel sets. Hence

$$\int_{\mathbb{R}^p} \mathcal{X}_E dm_p = \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} \mathcal{X}_E dm_1(x_{i_1}) \cdots dm_p(x_{i_p})$$

for any Borel set E after letting $n \rightarrow \infty$ and using the monotone convergence theorem. Approximating a nonnegative Borel function g with an increasing sequence of simple Borel measurable functions, and using the monotone convergence theorem yields

$$\int g dm_p = \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} g(x_1, \dots, x_p) dm_1(x_{i_1}) \cdots dm_p(x_{i_p})$$

The claim about the measure of a box being the product of the lengths of its sides also comes from this.

By Proposition 11.1.2, for f measurable, there exists g Borel measurable such that $g = f$ a.e. and $g \leq f$. Then

$$\int_{\mathbb{R}^p} f dm_p = \int_{\mathbb{R}^p} g dm_p = \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} g(x_1, \dots, x_p) dm_1(x_{i_1}) \cdots dm_p(x_{i_p})$$

It is similar if $h \geq f$ and equal to f a.e.

It remains to consider the claim about translation invariance. If R is a box, $R = \prod_{i=1}^p (a_i, b_i)$, then it is clear that $m_p(x + R) = m_p(R)$. Let \mathcal{H} be as above and let \mathcal{G} be those Borel sets E for which $m_p(x + E \cap R_n) = m_p(E \cap R_n)$ where R_n is as above. Thus \mathcal{G}

contains \mathcal{H} . Then it is obvious \mathcal{G} is closed with respect to countable disjoint unions. The case of complements maybe is not as obvious.

$$(\mathbf{x} + R_n) \setminus (\mathbf{x} + R_n \cap E) = \mathbf{x} + R_n \cap E^C.$$

Then

$$\begin{aligned} m_p(\mathbf{x} + R_n \cap E^C) &= m_p(\mathbf{x} + R_n) - m_p(\mathbf{x} + R_n \cap E) \\ &= m_p(R_n) - m_p(R_n \cap E) = m_p(R_n \cap E^C) \end{aligned}$$

Thus by Dynkin's lemma, $\mathcal{G} = \mathcal{B}(\mathbb{R}^p)$. Thus for all E Borel,

$$m_p(E \cap R_n) = m_p(\mathbf{x} + E \cap R_n).$$

Now let $n \rightarrow \infty$. It follows that m_p is translation invariant for all Borel sets.

In general, if E is Lebesgue measurable, it follows from Proposition 11.1.2 that there are sets $F \subseteq E \subseteq G$ where F, G are Borel and $m_p(F) = m_p(E) = m_p(G)$. Then

$$m_p(E) = m_p(F) = m_p(F + \mathbf{x}) \leq m_p(E + \mathbf{x}) \leq m_p(G + \mathbf{x}) = m_p(G) = m_p(E)$$

and so all the inequalities are equal signs. Hence $m_p(E + \mathbf{x}) = m_p(E)$. ■

The following is a useful lemma. In this lemma, X_i will be some metric space or more generally a topological space. It is useful in recognizing a Borel measurable set when you see it.

Lemma 11.3.2 *If E_i is a Borel set in X_i , then $\prod_{k=1}^p E_{i_k}$ is a Borel set in $\prod_{k=1}^p X_{i_k}$.*

Proof: Let $\pi_{i_r} : \prod_{k=1}^p X_{i_k} \rightarrow X_{i_r}$ be the projection map. That is $\pi_{i_r}(\mathbf{x}) = x_{i_r}$ when $\mathbf{x} = (x_{i_1}, x_{i_2}, \dots, x_{i_p})$. Obviously this is continuous. Therefore, if U is an open set in X_{i_r} , $\pi_{i_r}^{-1}(U) = X_{i_1} \times X_{i_2} \times \dots \times U \times \dots \times X_{i_p}$. Is an open set. Let \mathcal{B}_{i_r} be the Borel sets of X_{i_r} . E such that $\pi_{i_r}^{-1}(E) = X_{i_1} \times X_{i_2} \times \dots \times E \times \dots \times X_{i_p}$ is a Borel set in $\prod_{k=1}^p X_{i_k}$. Then \mathcal{B}_{i_r} is a σ algebra and it contains the open sets. Therefore, it contains the Borel sets of X_{i_r} . It follows that $\prod_{k=1}^p E_{i_k} = \cap_{k=1}^p \pi_{i_k}^{-1}(E_{i_k})$ is a finite intersection of Borel sets in $\prod_{k=1}^p X_{i_k}$ and so it is also a Borel set. ■

Example 11.3.3 *Let $A \equiv \{(x, y) : y < g(x)\}$ for g a Borel measurable real valued function. Then A is a Borel set.*

To see this, partition \mathbb{R} into equally spaced points $\{r_k^n\}_{k=-\infty}^{\infty}$, $r_k^n < r_{k+1}^n$, $r_{k+1}^n - r_k^n = 2^{-n}$ and let $g_n(x) \equiv \sum_{k=-\infty}^{\infty} r_{k-1}^n \mathcal{I}_{g^{-1}((r_{k-1}^n, r_k^n])}(x)$ so that $g_n(x) \rightarrow g(x)$ for each x . Let $A_n \equiv \{(x, y) : y < g_n(x)\}$. Now each A_k is Borel by the above Lemma. Then thanks to convergence, $A = \cap_{m=1}^{\infty} \cap_{k \geq m} A_k$ so A is Borel.

Example 11.3.4 *Let $A \equiv \{(x, y) : y \leq g(x)\}$ for g a Borel measurable function. Then A is a Borel set.*

This follows from observing that if $A_n \equiv \{(x, y) : y < g(x) + 2^{-n}\}$ then $A = \cap_{n=1}^{\infty} A_n$. Thus sets of the form $[a, b] \times \{(x, y) : y \leq g(x)\}$ for g Borel measurable are Borel measurable. These examples justify the usual calculus manipulations involving iterated integrals, the next example being an illustration of this.

Example 11.3.5 Find the iterated integral $\int_0^1 \int_x^1 \frac{\sin(y)}{y} dy dx$.

Notice the limits. The iterated integral equals $\int_{\mathbb{R}^2} \mathcal{X}_A(x, y) \frac{\sin(y)}{y} dm_2$ where

$$A = \{(x, y) : x \leq y \text{ where } x \in [0, 1]\}.$$

Fubini's theorem can be applied because the function $(x, y) \rightarrow \sin(y)/y$ is continuous except at $y = 0$ and can be redefined to be continuous there. The function is also bounded so $(x, y) \rightarrow \mathcal{X}_A(x, y) \frac{\sin(y)}{y}$ clearly is in $L^1(\mathbb{R}^2)$. Therefore,

$$\begin{aligned} \int_{\mathbb{R}^2} \mathcal{X}_A(x, y) \frac{\sin(y)}{y} dm_2 &= \int \int \mathcal{X}_A(x, y) \frac{\sin(y)}{y} dx dy \\ &= \int_0^1 \int_0^y \frac{\sin(y)}{y} dx dy = \int_0^1 \sin(y) dy = 1 - \cos(1) \end{aligned}$$

Here is a general and important result.

Example 11.3.6 *Integration by parts.* Suppose f, g are both absolutely continuous, then $\int_a^b f g' dt =$

$$\begin{aligned} \int_a^b \left(f(a) + \int_a^t f'(s) ds \right) g'(t) dt &= f(a)(g(b) - g(a)) + \int_a^b \int_a^t f'(s) ds g'(t) dt \\ &= f(a)(g(b) - g(a)) + \int_a^b \int_s^b f'(s) g'(t) dt ds \\ &= f(a)(g(b) - g(a)) + \int_a^b f'(s)(g(b) - g(s)) ds \\ &= f(a)(g(b) - g(a)) + g(b)(f(b) - f(a)) - \int_a^b f'(s) g(s) ds \\ &= g(b)f(b) - f(a)g(a) - \int_a^b f'(s) g(s) ds \end{aligned}$$

11.4 Maximal Functions

In this section the Besicovitch covering theorem, Theorem 4.5.8 will be used to obtain the Lebesgue differentiation theorem for general Radon measures. This will end up including Lebesgue measure presented later. In what follows, μ will be a Radon measure, complete, inner and outer regular, and finite on compact sets. Also

$$Z \equiv \{x \in \mathbb{R}^p : \mu(B(x, r)) = 0 \text{ for some } r > 0\}, \quad (11.4)$$

Lemma 11.4.1 Z is measurable and $\mu(Z) = 0$.

Proof: For each $x \in Z$, there exists a ball $B(x, r)$ with $\mu(B(x, r)) = 0$. Let \mathcal{C} be the collection of these balls. Since \mathbb{R}^p has a countable basis, a countable subset $\tilde{\mathcal{C}} \equiv \{B_i\}_{i=1}^\infty$, of \mathcal{C} also covers Z . Then letting $\bar{\mu}$ denote the outer measure determined by μ , $\bar{\mu}(Z) \leq \sum_{i=1}^\infty \bar{\mu}(B_i) = \sum_{i=1}^\infty \mu(B_i) = 0$. Therefore, Z is measurable, $(\bar{\mu}(S) \geq \bar{\mu}(S \cap Z) + \bar{\mu}(S \cap Z^C))$ and has measure zero as claimed. ■

Let $Mf : \mathbb{R}^p \rightarrow [0, \infty]$ by

$$Mf(x) \equiv \begin{cases} \sup_{r \leq 1} \frac{1}{\mu(B(x, r))} \int_{B(x, r)} |f| d\mu & \text{if } x \notin Z \\ 0 & \text{if } x \in Z \end{cases}.$$

I will begin using $\|f\|_1$ for the integral $\int_{\Omega} |f| d\mu$.

The special points described in the following theorem are called Lebesgue points.

Theorem 11.4.2 *Let μ be a Radon measure and let $f \in L^1(\mathbb{R}^p, \mu)$ meaning that $\int_{\Omega} |f| d\mu < \infty$. Then for μ a.e. x , $\lim_{r \rightarrow 0} \frac{1}{\mu(B(x, r))} \int_{B(x, r)} |f(y) - f(x)| d\mu(y) = 0$. Also $\bar{\mu}([Mf > \varepsilon]) \leq N_p \varepsilon^{-1} \|f\|_1$.*

Proof: First consider the following claim which is called a weak type estimate.

Claim 1: The following inequality holds for N_p the constant of the Besicovitch covering theorem, Theorem 4.5.8: $\bar{\mu}([Mf > \varepsilon]) \leq N_p \varepsilon^{-1} \|f\|_1$

Proof of claim: First note $[Mf > \varepsilon] \cap Z = \emptyset$ and without loss of generality, you can assume $\bar{\mu}([Mf > \varepsilon]) > 0$. Let U be an open set containing $[Mf > \varepsilon]$ such that $\bar{\mu}([Mf > \varepsilon])$. Next, for each $x \in [Mf > \varepsilon]$ there exists a ball $B_x = B(x, r_x)$ with $r_x \leq 1$ and the following inequality holding. $\mu(B_x)^{-1} \int_{B(x, r_x)} |f| d\mu > \varepsilon$. Let \mathcal{F} be this collection of balls so that $[Mf > \varepsilon]$ is the set of centers of balls of \mathcal{F} . By the Besicovitch covering theorem, Theorem 4.5.8, $[Mf > \varepsilon] \subseteq \cup_{i=1}^{N_p} \{B : B \in \mathcal{G}_i\}$ where \mathcal{G}_i is a collection of disjoint balls of \mathcal{F} . Now for some i , $\bar{\mu}([Mf > \varepsilon]) / N_p \leq \mu(\cup \{B : B \in \mathcal{G}_i\})$ because if this is not so, then for all i , $\bar{\mu}([Mf > \varepsilon]) / N_p > \mu(\cup \{B : B \in \mathcal{G}_i\})$ and so

$$\bar{\mu}([Mf > \varepsilon]) \leq \sum_{i=1}^{N_p} \mu(\cup \{B : B \in \mathcal{G}_i\}) < \sum_{i=1}^{N_p} \frac{\bar{\mu}([Mf > \varepsilon])}{N_p} = \bar{\mu}([Mf > \varepsilon]),$$

a contradiction. Therefore for this i ,

$$\begin{aligned} \frac{\bar{\mu}([Mf > \varepsilon])}{N_p} &\leq \mu(\cup \{B : B \in \mathcal{G}_i\}) = \sum_{B \in \mathcal{G}_i} \mu(B) \leq \sum_{B \in \mathcal{G}_i} \varepsilon^{-1} \int_B |f| d\mu \\ &\leq \varepsilon^{-1} \int_{\mathbb{R}^p} |f| d\mu = \varepsilon^{-1} \|f\|_1. \end{aligned}$$

This shows Claim 1.

Claim 2: If g is any continuous function defined on \mathbb{R}^p , then for $x \notin Z$,

$$\lim_{r \rightarrow 0} \frac{1}{\mu(B(x, r))} \int_{B(x, r)} |g(y) - g(x)| d\mu(y) = 0$$

and

$$\lim_{r \rightarrow 0} \frac{1}{\mu(B(x, r))} \int_{B(x, r)} g(y) d\mu(y) = g(x). \quad (11.5)$$

Proof: Since g is continuous at x , whenever r is small enough,

$$\frac{1}{\mu(B(x, r))} \int_{B(x, r)} |g(y) - g(x)| d\mu(y) \leq \frac{1}{\mu(B(x, r))} \int_{B(x, r)} \varepsilon d\mu(y) = \varepsilon.$$

11.5 follows from the above and the triangle inequality. This proves the claim.

Now let $g \in C_c(\mathbb{R}^p)$ and $x \notin Z$. Then from the above observations about continuous functions,

$$\bar{\mu} \left(\left[x \notin Z : \limsup_{r \rightarrow 0} \frac{1}{\mu(B(x, r))} \int_{B(x, r)} |f(y) - f(x)| d\mu(y) > \varepsilon \right] \right) \quad (11.6)$$

$$\begin{aligned} &\leq \bar{\mu} \left(\left[x \notin Z : \limsup_{r \rightarrow 0} \frac{1}{\mu(B(x, r))} \int_{B(x, r)} |f(y) - g(y)| d\mu(y) > \frac{\varepsilon}{2} \right] \right) \\ &\quad + \bar{\mu} \left(\left[x \notin Z : |g(x) - f(x)| > \frac{\varepsilon}{2} \right] \right). \\ &\leq \bar{\mu} \left(\left[M(f - g) > \frac{\varepsilon}{2} \right] \right) + \bar{\mu} \left(\left[|f - g| > \frac{\varepsilon}{2} \right] \right) \end{aligned} \quad (11.7)$$

Now $\int_{[|f-g| > \frac{\varepsilon}{2}]} |f - g| d\mu \geq \frac{\varepsilon}{2} \bar{\mu}([|f - g| > \frac{\varepsilon}{2}])$ and so using Claim 1 in 11.7, it follows that 11.6 is dominated by $\left(\frac{2}{\varepsilon} + \frac{N_p}{\varepsilon}\right) \int |f - g| d\mu$. But by Theorem 10.8.7, g can be chosen to make this as small as desired. Hence 11.6 is 0. Now observe that

$$\begin{aligned} &\bar{\mu} \left(\left[x \notin Z : \limsup_{r \rightarrow 0} \frac{1}{\mu(B(x, r))} \int_{B(x, r)} |f(y) - f(x)| d\mu(y) > 0 \right] \right) \\ &\leq \sum_{k=1}^{\infty} \bar{\mu} \left(\left[x \notin Z : \limsup_{r \rightarrow 0} \frac{1}{\mu(B(x, r))} \int_{B(x, r)} |f(y) - f(x)| d\mu(y) > \frac{1}{k} \right] \right) = 0 \end{aligned}$$

By completeness of μ this implies

$$\left[x \notin Z : \limsup_{r \rightarrow 0} \frac{1}{\mu(B(x, r))} \int_{B(x, r)} |f(y) - f(x)| d\mu(y) > 0 \right]$$

is a set of μ measure zero. ■

The following corollary is the main result referred to as the Lebesgue Besicovitch Differentiation theorem.

Definition 11.4.3 $f \in L^1_{loc}(\mathbb{R}^p, \mu)$ means $f \mathcal{X}_B$ is in $L^1(\mathbb{R}^p, \mu)$ whenever B is a ball.

Theorem 11.4.4 If $f \in L^1_{loc}(\mathbb{R}^p, \mu)$, then for μ a.e. $x \notin Z$,

$$\lim_{r \rightarrow 0} \frac{1}{\mu(B(x, r))} \int_{B(x, r)} |f(y) - f(x)| d\mu(y) = 0. \quad (11.8)$$

Proof: If f is replaced by $f \mathcal{X}_{B(0, k)}$ then the conclusion 11.8 holds for all $x \notin F_k$ where F_k is a set of μ measure 0. Letting $k = 1, 2, \dots$, and $F \equiv \bigcup_{k=1}^{\infty} F_k$, it follows that F is a set of measure zero and for any $x \notin F$, and $k \in \{1, 2, \dots\}$, 11.8 holds if f is replaced by $f \mathcal{X}_{B(0, k)}$. Picking any such x , and letting $k > |x| + 1$, this shows

$$\begin{aligned} &\lim_{r \rightarrow 0} \frac{1}{\mu(B(x, r))} \int_{B(x, r)} |f(y) - f(x)| d\mu(y) \\ &= \lim_{r \rightarrow 0} \frac{1}{\mu(B(x, r))} \int_{B(x, r)} |f \mathcal{X}_{B(0, k)}(y) - f \mathcal{X}_{B(0, k)}(x)| d\mu(y) = 0 \end{aligned}$$

because for all r small enough, $B(x, r) \subseteq B(0, k)$. ■

Definition 11.4.5 Let E be a measurable set. Then $\mathbf{x} \in E$ is called a point of density if $\mathbf{x} \notin Z$ and $\lim_{r \rightarrow 0} \frac{\mu(B(\mathbf{x}, r) \cap E)}{\mu(B(\mathbf{x}, r))} = 1$. Recall Z is the set of points \mathbf{x} where $\mu(B(\mathbf{x}, r)) = 0$ for some $r > 0$.

Proposition 11.4.6 Let E be a measurable set. Then μ a.e. $\mathbf{x} \in E$ is a point of density.

Proof: This follows from letting $f(\mathbf{x}) = \chi_E(\mathbf{x})$ in Theorem 11.4.4. ■

From Theorem 11.4.2, $\mu([Mf > \lambda]) \leq \frac{N_p}{\lambda} \|f\|_{L^1}$.

$$\begin{aligned} Mf &\equiv \sup_{r \leq 1} \frac{1}{\mu(B(\mathbf{x}, r))} \int_{B(\mathbf{x}, r)} |f| d\mu \leq \sup_{r \leq 1} \frac{1}{\mu(B(\mathbf{x}, r))} \int_{B(\mathbf{x}, r)} \left(|f| \chi_{[|f| > \frac{\lambda}{2}]} + \frac{\lambda}{2} \right) d\mu \\ &= M\left(f \chi_{[|f| > \frac{\lambda}{2}]}\right) + \frac{\lambda}{2} \end{aligned}$$

Therefore, $[Mf > \lambda] \subseteq \left[M\left(f \chi_{[|f| > \frac{\lambda}{2}]}\right) > \lambda/2 \right]$ so

$$[Mf > 2\lambda] \subseteq [M(f \chi_{[|f| > \lambda]}) > \lambda] \leq \frac{N_p}{\lambda} \int_{[|f| > \lambda]} |f| d\mu.$$

This shows the following modified weak estimate.

Corollary 11.4.7 Let f be in $L^1(\mathbb{R}^p, \mu)$. Then $\mu([Mf > 2\lambda]) \leq \frac{N_p}{\lambda} \int_{[|f| > \lambda]} |f| d\mu$.

11.5 Strong Estimates for Maximal Function

Here $p > 1$, not the dimension. Let $\|f\|_{L^p(\mathbb{R}^n)}^p \equiv \int_{\mathbb{R}^n} |f|^p d\mu$. Let $\lambda^{1/p} \equiv 2\eta$ so $\lambda = 2^p \eta^p$ and $d\lambda = 2^p p \eta^{p-1} d\eta$. Then use Corollary 11.4.7 so

$$\begin{aligned} \int |Mf|^p d\mu &= \int_0^\infty \mu\left([Mf > \lambda^{1/p}]\right) d\lambda = \int_0^\infty \mu([Mf > 2\eta]) 2^p p \eta^{p-1} d\eta \\ &\leq \int_0^\infty \frac{N}{\eta} \int_{[|f| > \eta]} |f| 2^p p \eta^{p-1} d\mu d\eta = N 2^p p \int \int_0^{|f|} |f| \eta^{p-2} d\eta d\mu = C_p \int |f|^p d\mu \end{aligned}$$

Of course this is all assuming that Mf is measurable. This is most easily shown if $\mu = m_n$ Lebesgue measure and this is the case of most interest to me. Consider $\mathbf{x} \rightarrow \int_{B(\mathbf{x}, r)} |f(\mathbf{y})| dm \equiv f_r(\mathbf{x})$. This is continuous assuming $f \in L^1_{loc}$ thanks to continuity of translation of Lebesgue measure. Thus $\mathbf{x} \rightarrow \frac{1}{\mu(B(\mathbf{x}, r))} \int_{B(\mathbf{x}, r)} |f(\mathbf{y})| dm$ is continuous. Now it follows that $Mf(\mathbf{x}) = \sup_{0 < r < 1} f_r(\mathbf{x}) = \sup_{0 < r < 1, r \in \mathbb{Q}} f_r(\mathbf{x})$ is Borel measurable. Therefore, we can state the following corollary.

Corollary 11.5.1 Let $\|f\|_{L^p} < \infty$ where $\mu = m_n$ in \mathbb{R}^n . Then $\|Mf\|_{L^p} \leq C_p \|f\|_{L^p}$ where C depends only on $p > 1$ and the dimension. This is called a strong estimate for the maximal function as opposed to the one from Theorem 11.4.2 for $p = 1$ which is called a weak estimate.

11.6 The Brouwer Fixed Point Theorem

I found this proof of the Brouwer fixed point theorem in Evans [18] and Dunford and Schwartz [16]. The main idea which makes proofs like this work is Lemma 7.11.2 which is stated next for convenience.

Lemma 11.6.1 *Let $g : U \rightarrow \mathbb{R}^p$ be C^2 where U is an open subset of \mathbb{R}^p . Then it follows that $\sum_{j=1}^p \text{cof}(Dg)_{ij,j} = 0$, where here $(Dg)_{ij} \equiv g_{i,j} \equiv \frac{\partial g_i}{\partial x_j}$. Also, $\text{cof}(Dg)_{ij} = \frac{\partial \det(Dg)}{\partial g_{i,j}}$.*

Definition 11.6.2 *Let h be a function defined on an open set, $U \subseteq \mathbb{R}^p$. Then $h \in C^k(\overline{U})$ if there exists a function g defined on an open set, W containing \overline{U} such that $g = h$ on U and g is $C^k(W)$.*

Lemma 11.6.3 *There does not exist $h \in C^2(\overline{B(0,R)})$ with $h : \overline{B(0,R)} \rightarrow \partial B(0,R)$ which has the property that $h(x) = x$ for all $x \in \partial B(0,R) \equiv \{x : |x| = R\}$. Such a function is called a retract.*

Proof: First note that if h is such a retract, then for all $x \in B(0,R)$, $\det(Dh(x)) = 0$. This is because if $\det(Dh(x)) \neq 0$ for some such x , then by the inverse function theorem, $h(B(x,\delta))$ is an open set for small enough δ but this would require that this open set is a subset of $\partial B(0,R)$ which is impossible because no open ball is contained in $\partial B(0,R)$. Here and below, let B_R denote $\overline{B(0,R)}$.

Now suppose such an h exists. Let $\lambda \in [0, 1]$ and let $p_\lambda(x) \equiv x + \lambda(h(x) - x)$. This function, p_λ is a homotopy of the identity map and the retract h . Let

$$I(\lambda) \equiv \int_{B(0,R)} \det(Dp_\lambda(x)) dx.$$

Then using the dominated convergence theorem,

$$\begin{aligned} I'(\lambda) &= \int_{B(0,R)} \sum_{i,j} \frac{\partial \det(Dp_\lambda(x))}{\partial p_{\lambda i,j}} \frac{\partial p_{\lambda i,j}(x)}{\partial \lambda} dx \\ &= \int_{B(0,R)} \sum_i \sum_j \frac{\partial \det(Dp_\lambda(x))}{\partial p_{\lambda i,j}} (h_i(x) - x_i)_j dx \\ &= \int_{B(0,R)} \sum_i \sum_j \text{cof}(Dp_\lambda(x))_{ij} (h_i(x) - x_i)_j dx \end{aligned}$$

Now by assumption, $h_i(x) = x_i$ on $\partial B(0,R)$ and so one can integrate by parts, in the iterated integrals used to compute $\int_{B(0,R)}$ and write

$$I'(\lambda) = - \sum_i \int_{B(0,R)} \sum_j \text{cof}(Dp_\lambda(x))_{ij,j} (h_i(x) - x_i) dx = 0.$$

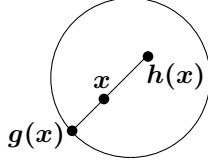
Therefore, $I(\lambda)$ equals a constant. However, $I(0) = m_p(B(0,R)) \neq 0$ and as pointed out above, $I(1) = 0$. ■

The following is the Brouwer fixed point theorem for C^2 maps.

Lemma 11.6.4 *If $h \in C^2(\overline{B(0, R)})$ and $h : \overline{B(0, R)} \rightarrow \overline{B(0, R)}$, then h has a fixed point x such that $h(x) = x$.*

Proof: Suppose the lemma is not true. Then for all x , $|x - h(x)| \neq 0$. Then define $g(x) = h(x) + \frac{x - h(x)}{|x - h(x)|} t(x)$ where $t(x)$ is nonnegative and is chosen such that $g(x) \in \partial B(0, R)$.

This mapping is illustrated in the following picture.



If $x \rightarrow t(x)$ is C^2 near $\overline{B(0, R)}$, it will follow g is a C^2 retract onto $\partial B(0, R)$ contrary to Lemma 11.6.3. Thus $t(x)$ is the nonnegative solution t to

$$\left| h(x) + \frac{x - h(x)}{|x - h(x)|} t(x) \right|^2 = |h(x)|^2 + 2 \left(h(x), \frac{x - h(x)}{|x - h(x)|} \right) t + t^2 = R^2 \quad (11.9)$$

then by the quadratic formula,

$$t(x) = - \left(h(x), \frac{x - h(x)}{|x - h(x)|} \right) + \sqrt{\left(h(x), \frac{x - h(x)}{|x - h(x)|} \right)^2 + (R^2 - |h(x)|^2)}$$

Is $x \rightarrow t(x)$ a function in C^2 ? If what is under the radical is positive, then this is so because $s \rightarrow \sqrt{s}$ is smooth for $s > 0$. In fact, this is the case here. The inside of the radical is positive if $R > |h(x)|$. If $|h(x)| = R$, it is still positive because in this case, the angle between $h(x)$ and $x - h(x)$ cannot be $\pi/2$ because of the shape of the ball. This shows that $x \rightarrow t(x)$ is the composition of C^2 functions and is therefore C^2 . Thus this $g(x)$ is a C^2 retract and by the above lemma, there isn't one. ■

Now it is easy to prove the Brouwer fixed point theorem. The following theorem is the Brouwer fixed point theorem for a ball.

Theorem 11.6.5 *Let B_R be the above closed ball and let $f : B_R \rightarrow B_R$ be continuous. Then there exists $x \in B_R$ such that $f(x) = x$.*

Proof: Let $f_k(x) \equiv \frac{f(x)}{1 + 1/k}$. Thus

$$\begin{aligned} \|f_k - f\| &= \max_{x \in B_R} \left\{ \left| \frac{f(x)}{1 + (1/k)} - f(x) \right| \right\} = \max_{x \in B_R} \left\{ \left| \frac{f(x) - f(x)(1 + (1/k))}{1 + (1/k)} \right| \right\} \\ &= \max_{x \in B_R} \left\{ \left| \frac{f(x)(1/k)}{1 + (1/k)} \right| \right\} \leq \frac{R}{1 + k} \end{aligned}$$

Letting $\|h\| \equiv \max \{|h(x)| : x \in B_R\}$, It follows from the Weierstrass approximation theorem, there exists a function whose components are polynomials g_k such that $\|g_k - f_k\| < \frac{R}{k+1}$. Then if $x \in B_R$, it follows

$$|g_k(x)| \leq |g_k(x) - f_k(x)| + |f_k(x)| < \frac{R}{1 + k} + \frac{kR}{1 + k} = R$$

and so g_k maps B_R to B_R . By Lemma 11.6.4 each of these g_k has a fixed point x_k such that $g_k(x_k) = x_k$. The sequence of points, $\{x_k\}$ is contained in the compact set, B_R and so there exists a convergent subsequence still denoted by $\{x_k\}$ which converges to a point $x \in B_R$. Then from uniform convergence of g_k to f , $f(x) = \lim_{k \rightarrow \infty} f(x_k) = \lim_{k \rightarrow \infty} g_k(x_k) = \lim_{k \rightarrow \infty} x_k = x$. ■

Definition 11.6.6 A set A is a retract of a set B if $A \subseteq B$, and there is a continuous map $h : B \rightarrow A$ such that $h(x) = x$ for all $x \in A$ and h is onto. B has the fixed point property means that whenever g is continuous and $g : B \rightarrow B$, it follows that g has a fixed point.

Proposition 11.6.7 Let A be a retract of B and suppose B has the fixed point property. Then so does A .

Proof: Suppose $f : A \rightarrow A$. Let h be the retract of B onto A . Then $f \circ h : B \rightarrow B$ is continuous. Thus, it has a fixed point $x \in B$ so $f(h(x)) = x$. However, $h(x) \in A$ and $f : A \rightarrow A$ so in fact, $x \in A$. Now $h(x) = x$ and so $f(x) = x$. ■

Recall that every convex compact subset K of \mathbb{R}^p is a retract of all of \mathbb{R}^p obtained by using the projection map. See Problems beginning with 8 on Page 151. In particular, K is a retract of a large closed ball containing K , which ball has the fixed point property. Therefore, K also has the fixed point property. This shows the following which is a convenient formulation of the Brouwer fixed point theorem. However, Proposition 11.6.7 is more general.

Theorem 11.6.8 Every convex closed and bounded subset of \mathbb{R}^p has the fixed point property.

11.7 Change of Variables, Linear Maps

This is about changing the variables for linear maps where \mathcal{F}_p denotes the Lebesgue measurable sets.

Theorem 11.7.1 In case $h : \mathbb{R}^p \rightarrow \mathbb{R}^p$ is Lipschitz, satisfying the Lipschitz condition $\|h(x) - h(y)\| \leq K\|x - y\|$, then if T is a set for which $m_p(T) = 0$, it follows that $m_p(h(T)) = 0$. Also if $E \in \mathcal{F}_p$, then $h(E) \in \mathcal{F}_p$.

Proof: By the Lipschitz condition, $\|h(x + v) - h(x)\| \leq K\|v\|$ and you can simply let $T \subseteq V$ where $m_p(V) < \varepsilon / (K^p 5^p)$. Then there is a countable disjoint sequence of balls $\{B_i\}$ such that $\{\hat{B}_i\}$ covers T and each ball B_i is contained in V each having radius no more than 1. Then the Lipschitz condition implies $h(\hat{B}_i) \subseteq B(h(x_i), 5K)$ and so

$$\bar{m}_p(h(T)) \leq \sum_{i=1}^{\infty} m_p(h(\hat{B}_i)) \leq 5^p K^p \sum_{i=1}^{\infty} m_p(B_i) \leq K^p 5^p m_p(V) < \varepsilon$$

Since ε is arbitrary, this shows that $h(T)$ is measurable and $m_p(h(T)) = 0$.

Now let $E \in \mathcal{F}_p$, $m_p(E) < \infty$. Then by of the measure and Theorem 9.8.6, there exists F which is the countable union of compact sets such that $E = F \cup N$ where N is a set of measure zero. Then from the first part, $h(E \setminus F) \subseteq h(N)$ and this set on the right has measure zero and so by completeness of the measure, $h(E \setminus F) \in \mathcal{F}_p$ and so $h(E) = h(E \setminus F) \cup h(F) \in \mathcal{F}_p$ because $F = \cup_k K_k$, each K_k compact. Hence $h(F) = \cup_k h(K_k)$

which is the countable union of compact sets, a Borel set, due to the continuity of h . For arbitrary E , $h(E) = \bigcup_{k=1}^{\infty} h(E \cap B(\mathbf{0}, k)) \in \mathcal{F}_p$. ■

Of course an example of a Lipschitz map is a linear map. $\|A\mathbf{x} - A\mathbf{y}\| = \|A(\mathbf{x} - \mathbf{y})\| \leq \|A\| \|\mathbf{x} - \mathbf{y}\|$. Therefore, if A is linear and E is Lebesgue measurable, then $A(E)$ is also Lebesgue measurable. This is convenient.

Lemma 11.7.2 *Every open set U in \mathbb{R}^p is a countable disjoint union of half open boxes of the form $Q \equiv \prod_{i=1}^p [a_i, a_i + 2^{-k})$ where $a_i = l2^{-k}$ for l some integer.*

Proof: It is clear that there exists \mathcal{Q}_k a countable disjoint collection of these half open boxes each of sides of length 2^{-k} whose union is all of \mathbb{R}^p . Let \mathcal{B}_1 be those sets of \mathcal{Q}_1 which are contained in U , if any. Having chosen \mathcal{B}_{k-1} , let \mathcal{B}_k consist of those sets of \mathcal{Q}_k which are contained in U such that none of these are contained in \mathcal{B}_{k-1} . Then $\bigcup_{k=1}^{\infty} \mathcal{B}_k$ is a countable collection of disjoint boxes of the right sort whose union is U . This is because if R is a box of \mathcal{Q}_k and \hat{R} is a box of \mathcal{Q}_{k-1} , then either $R \subseteq \hat{R}$ or $R \cap \hat{R} = \emptyset$. If $p \in U$ then it is ultimately contained in some \mathcal{B}_k for k as small as possible because p is at a positive distance from U^C . ■

Corollary 11.7.3 *If D is a diagonal matrix having nonnegative eigenvalues, and U is an open set, it follows that $m_p(DU) = \det(D)m_p(U)$.*

Proof: The multiplication by D just scales each side of the boxes whose disjoint union is U , multiplying the side in the i^{th} direction by the i^{th} diagonal element. Thus if R is one of the boxes, $m_p(DR) = \det(D)m_p(R)$. The desired result follows from adding these together. ■

I will write dx or dy instead of $dm_p(x)$ or $dm_p(y)$ to save on notation.

Theorem 11.7.4 *Let $E \in \mathcal{F}_p$ and let A be a $p \times p$ matrix. Then $A(E)$ is Lebesgue measurable and $m_p(A(E)) = |\det(A)|m_p(E)$. Also, if E is any Lebesgue measurable set, then $\int \mathcal{X}_{A(E)}(\mathbf{y}) d\mathbf{y} = \int \mathcal{X}_E(\mathbf{x}) |\det(A)| d\mathbf{x}$.*

Proof: First note that if $C(\mathbf{x}, r) \equiv \{\mathbf{y} \in \mathbb{R}^p : |\mathbf{y} - \mathbf{x}| = r\}$, then $m_p(C(\mathbf{x}, r)) = 0$. This follows from translation invariance and Corollary 11.7.3 applied to diagonal D having diagonal entries $r(1 + \varepsilon)$ and one with diagonal entries $r(1 - \varepsilon)$ to obtain that for arbitrary $\varepsilon > 0$,

$$\begin{aligned} m_p(C(\mathbf{0}, r)) &\leq m_p(B(\mathbf{0}, (1 + \varepsilon)r) \setminus B(\mathbf{0}, (1 - \varepsilon)r)) \\ &= m_p(B(\mathbf{0}, r)) [((1 + \varepsilon)r)^p - ((1 - \varepsilon)r)^p] \end{aligned}$$

Here $|\cdot|$ is the Euclidean norm so all orthogonal transformations acting on a ball centered at $\mathbf{0}$ leave the ball unchanged. Now let U be an open set, then by Theorem 9.12.2, there are disjoint open balls $\{B_i\}_{i=1}^{\infty}$ such that $U = (\bigcup_i B_i) \cup N$ where $m_n(N) = 0$.

From the right polar decomposition, Theorem 1.5.5 and the fact that one can diagonalize a symmetric matrix S , $A = RS = RQ^*DQ$ where R and Q are orthogonal matrices and D is a diagonal matrix with all nonnegative diagonal entries. Thus, if B is an open ball centered at $\mathbf{0}$,

$$\begin{aligned} m_p(A(B)) &= m_p(RQ^*DQ(B)) = m_p(RQ^*D(B)) \\ &= |\det(R)| |\det(Q^*)| \det(D) m_p(B) = |\det(A)| m_p(B) \end{aligned}$$

By continuity of translation, the same holds if B has a center at some other point than $\mathbf{0}$. It follows that $m_p(A(U)) = \sum_i m_p(AB_i) = \sum_i |\det(A)| m_p(B_i) = |\det(A)| m_p(U)$. Now let \mathcal{K} be the open sets and \mathfrak{S} be those Borel sets E such that $m_p(A(E \cap B(\mathbf{0}, n))) = |\det(A)| m_p(E \cap B(\mathbf{0}, n))$. It is routine to verify that \mathfrak{S} is closed with respect to countable disjoint unions and complements. Therefore, $\mathfrak{S} = \sigma(\mathcal{K})$ and so this holds for all Borel E . Letting $n \rightarrow \infty$, it follows that for all E Borel, $m_p(A(E)) = |\det(A)| m_p(E)$.

If E is only Lebesgue measurable, then by regularity and Proposition 11.1.2, there exists G and F , G_δ and F_σ sets respectively such that $F \subseteq E \subseteq G$ and $m_p(G) = m_p(E) = m_p(F)$. Then $AF \subseteq AE \subseteq AG$ and for \bar{m}_p the outer measure determined by m_p ,

$$\begin{aligned} |\det(A)| m_p(F) &= m_p(AF) \leq m_p(AE) \leq m_p(AG) \\ &= \det(A) m_p(G) = |\det(A)| m_p(E) = |\det(A)| m_p(F) \end{aligned}$$

Thus all inequalities are equal signs. ■

Theorem 11.7.5 *Let $f \geq 0$ and suppose it is Lebesgue measurable. Then if A is a $p \times p$ matrix,*

$$\int \mathcal{X}_{A(\mathbb{R}^p)}(\mathbf{y}) f(\mathbf{y}) dm_p(\mathbf{y}) = \int f(A\mathbf{x}) |\det(A)| dm_p(\mathbf{x}). \quad (11.10)$$

Proof: From Theorem 11.7.4, the equation is true if $\det(A) = 0$. It follows that it suffices to consider only the case where A^{-1} exists. First suppose $f(\mathbf{y}) = \mathcal{X}_E(\mathbf{y})$ where E is a Lebesgue measurable set. In this case, $A(\mathbb{R}^n) = \mathbb{R}^n$. Then from Theorem 11.7.4

$$\begin{aligned} \int \mathcal{X}_{A(\mathbb{R}^p)}(\mathbf{y}) f(\mathbf{y}) d\mathbf{y} &= \int \mathcal{X}_E(\mathbf{y}) d\mathbf{y} = m_p(E) = |\det(A)| m_p(A^{-1}E) \\ &= \int_{\mathbb{R}^n} |\det(A)| \mathcal{X}_{A^{-1}E}(\mathbf{x}) d\mathbf{x} = \int_{\mathbb{R}^n} |\det(A)| \mathcal{X}_E(A\mathbf{x}) d\mathbf{x} = \int f(A\mathbf{x}) |\det(A)| d\mathbf{x} \end{aligned}$$

It follows from this that 11.10 holds whenever f is a nonnegative simple function. Finally, the general result follows from approximating the Lebesgue measurable function with nonnegative simple functions using Theorem 9.1.6 and then applying the monotone convergence theorem. ■

This is now a very good change of variables formula for a linear transformation. Next this is extended to differentiable functions.

11.8 Differentiable Functions and Measurability

To begin with, certain kinds of functions map measurable sets to measurable sets. It was shown earlier, Theorem 11.7.1, that Lipschitz functions do this. So do differentiable functions.

In this part of the argument it is convenient to take all balls with respect to the norm on \mathbb{R}^p given by $\|\mathbf{x}\| = \max\{|x_k| : k = 1, 2, \dots, p\}$. Thus from the definition of this norm, $B(\mathbf{x}, r)$ is the open box, $\prod_{k=1}^p (x_k - r, x_k + r)$ and so $m_p(B(\mathbf{x}, r)) = (2r)^p = 2^p r^p$. Also for a linear transformation $A \in \mathcal{L}(\mathbb{R}^p, \mathbb{R}^p)$, I will continue to use $\|A\| \equiv \sup_{\|\mathbf{x}\| \leq 1} \|A\mathbf{x}\|$.

Lemma 11.8.1 *Let $T \subseteq U$, where U is open, \mathbf{h} is continuous, and let \mathbf{h} be differentiable at each $\mathbf{x} \in T$ and suppose that $m_p(T) = 0$, then $m_p(\mathbf{h}(T)) = 0$.*

Proof: For $k \in \mathbb{N}$, let $T_k \equiv \{x \in T : \|Dh(x)\| < k\}$ and let $\varepsilon > 0$ be given. Since T_k is a subset of a set of measure zero, it is measurable, but we don't need to pay much attention to this fact. Now by outer regularity, there exists an open set V , containing T_k which is contained in U such that $m_p(V) < \varepsilon$. Let $x \in T_k$. Then by differentiability, $h(x+v) = h(x) + Dh(x)v + o(v)$ and so there exist arbitrarily small $r_x < 1$ such that $B(x, 5r_x) \subseteq V$ and whenever $\|v\| \leq 5r_x$, $\|o(v)\| < \frac{1}{5}\|v\|$. Thus, from the Vitali covering theorem, Theorem 4.5.3,

$$\begin{aligned} h(B(x, 5r_x)) &\subseteq Dh(x)(B(0, 5r_x)) + h(x) + B(0, r_x) \subseteq B(0, k5r_x) + \\ &+ B(0, r_x) + h(x) \subseteq B(h(x), (5k+1)r_x) \subseteq B(h(x), 6kr_x) \end{aligned}$$

From the Vitali covering theorem, there exists a countable disjoint sequence of these balls, $\{B(x_i, r_i)\}_{i=1}^\infty$ such that $\{B(x_i, 5r_i)\}_{i=1}^\infty = \{\widehat{B}_i\}_{i=1}^\infty$ covers T_k . Then letting \overline{m}_p denote the outer measure determined by m_p ,

$$\begin{aligned} \overline{m}_p(h(T_k)) &\leq \overline{m}_p\left(h\left(\bigcup_{i=1}^\infty \widehat{B}_i\right)\right) \leq \sum_{i=1}^\infty \overline{m}_p\left(h(\widehat{B}_i)\right) \\ &\leq \sum_{i=1}^\infty m_p(B(h(x_i), 6kr_{x_i})) = \sum_{i=1}^\infty m_p(B(x_i, 6kr_{x_i})) \\ &= (6k)^p \sum_{i=1}^\infty m_p(B(x_i, r_{x_i})) \leq (6k)^p m_p(V) \leq (6k)^p \varepsilon. \end{aligned}$$

Since $\varepsilon > 0$ is arbitrary, this shows $m_p(h(T_k)) = \overline{m}_p(h(T_k)) = 0$. Now $m_p(h(T)) = \lim_{k \rightarrow \infty} m_p(h(T_k)) = 0$. ■

Lemma 11.8.2 *Let h be continuous on U and let h be differentiable on $T \subseteq U$. If S is a Lebesgue measurable subset of T , then $h(S)$ is Lebesgue measurable.*

Proof: By Theorem 11.2.2 there exists F which is a countable union of compact sets $F = \bigcup_{k=1}^\infty K_k$ such that $F \subseteq S$, $m_p(S \setminus F) = 0$. Then $h(F) = \bigcup_k h(K_k) \in \mathcal{B}(\mathbb{R}^p)$ because the continuous image of a compact set is compact. Also, $h(S \setminus F)$ is a set of measure zero by Lemma 11.8.1 and so $h(S) = h(F) \cup h(S \setminus F) \in \mathcal{F}_p$ because it is the union of two sets which are in \mathcal{F}_p . ■

In particular, this proves the following theorem from a different point of view to that done before, using $x \rightarrow Ax$ being differentiable rather than $x \rightarrow Ax$ being Lipschitz. Later on, is a theorem which says that Lipschitz implies differentiable a.e. However, it is also good to note that if h has a derivative on an open set U , it does not follow that h is Lipschitz.

I will also use the following fundamental assertion, Sard's lemma.

Lemma 11.8.3 (Sard) *Let U be an open set in \mathbb{R}^p . Let $h : U \rightarrow \mathbb{R}^p$ be continuous and let h be differentiable on $A \subseteq U$. Let $Z \equiv \{x \in A : \det Dh(x) = 0\}$. Then $m_p(h(Z)) = 0$.*

Proof: Suppose first that A is bounded. Let $\varepsilon > 0$ be given. Also let $V \supseteq Z$ with $V \subseteq U$ open, and $m_p(Z) + \varepsilon > m_p(V)$. Now let $x \in Z$. Then since h is differentiable at x , there exists $\delta_x > 0$ such that if $r < \delta_x$, then $B(x, r) \subseteq V$ and also,

$$h(B(x, r)) \subseteq h(x) + Dh(x)(B(0, r)) + B(0, r\eta), \quad \eta < 1.$$

Regard $D\mathbf{h}(\mathbf{x})$ as an $n \times n$ matrix, the matrix of the linear transformation $D\mathbf{h}(\mathbf{x})$ with respect to the usual coordinates. Since $\mathbf{x} \in Z$, it follows that there exists an invertible matrix M such that $MD\mathbf{h}(\mathbf{x})$ is in row reduced echelon form with a row of zeros on the bottom. Therefore, using Theorem 11.7.4 about taking out the determinant of a transformation,

$$\begin{aligned} m_p(\mathbf{h}(B(\mathbf{x}, r))) &= |\det(M^{-1})| m_p(M(\mathbf{h}(B(\mathbf{x}, r)))) \\ &\leq |\det(M^{-1})| m_p(M(D\mathbf{h}(\mathbf{x}))(B(\mathbf{0}, r)) + MB(\mathbf{0}, r\eta)) \\ &\leq |\det(M^{-1})| \alpha_{p-1} \|M(D\mathbf{h}(\mathbf{x}))\|^{p-1} (2r + 2\eta r)^{p-1} \|M\| 2r\eta \\ &\leq C(\|M\|, |\det(M^{-1})|, \|D\mathbf{h}(\mathbf{x})\|) 4^{p-1} r^p 2\eta \end{aligned}$$

Here α_n is the volume of the unit ball in \mathbb{R}^n . This is because $M(D\mathbf{h}(\mathbf{x}))(B(\mathbf{0}, r)) + MB(\mathbf{0}, r\eta)$ in the third line up is contained in a cylinder, the base in \mathbb{R}^{p-1} which has radius $\|M(D\mathbf{h}(\mathbf{x}))\| (2r + 2\eta r)$ and height $\|M\| 2r\eta$. Thus its measure is no more than $\int_{\mathbb{R}^{p-1}} \int_{-\|Mr\eta\|}^{\|Mr\eta\|} dx_p dm_{p-1}$. Then letting $\delta_{\mathbf{x}}$ be still smaller if necessary, corresponding to sufficiently small η ,

$$m_p(\mathbf{h}(B(\mathbf{x}, r))) \leq \varepsilon m_p(B(\mathbf{x}, r)).$$

The balls of this form constitute a Vitali cover of Z . Hence, by the covering theorem Corollary 9.12.5, there exists $\{B_i\}_{i=1}^{\infty}$, $B_i = B_i(\mathbf{x}_i, r_i)$, a collection of disjoint balls, each of which is contained in V , such that $m_p(\mathbf{h}(B_i)) \leq \varepsilon m_p(B_i)$ and $m_p(Z \setminus \cup_i B_i) = 0$. Hence from Lemma 11.8.1,

$$m_p(\mathbf{h}(Z) \setminus \cup_i \mathbf{h}(B_i)) \leq m_p(\mathbf{h}(Z \setminus \cup_i B_i)) = 0$$

Therefore,

$$m_p(\mathbf{h}(Z)) \leq \sum_i m_p(\mathbf{h}(B_i)) \leq \varepsilon \sum_i m_p(B_i) \leq \varepsilon (m_p(V)) \leq \varepsilon (m_p(Z) + \varepsilon).$$

Since ε is arbitrary, this shows $m_p(\mathbf{h}(Z)) = 0$. What if A is not bounded? Then consider $Z_n = Z \cap B(\mathbf{0}, n) \subseteq A \cap B(\mathbf{0}, n)$. From what was just shown, $\mathbf{h}(Z_n)$ has measure 0 and so it follows that $\mathbf{h}(Z)$ also does, being the countable union of sets of measure zero. ■

11.9 Change of Variables, Nonlinear Maps

This preparation leads to a good change of variables formula. First is a lemma which is likely familiar by now.

Lemma 11.9.1 *Let $\mathbf{h} : \Omega \rightarrow \mathbb{R}^p$ where (Ω, \mathcal{F}) is a measurable space and suppose \mathbf{h} is continuous. Then $\mathbf{h}^{-1}(B) \in \mathcal{F}$ whenever B is a Borel set.*

Proof: Measurability applied to components of \mathbf{h} shows that $\mathbf{h}^{-1}(U) \in \mathcal{F}$ whenever U is an open set. If \mathcal{G} consists of the subsets G of \mathbb{R}^p for which $\mathbf{h}^{-1}(G) \in \mathcal{F}$, then \mathcal{G} is a σ algebra and \mathcal{G} contains the open sets. ■

Definition 11.9.2 *Let $\mathbf{h} : U \rightarrow \mathbb{R}^p$ be continuous, U open, and let $H \subseteq U$ be measurable and \mathbf{h} is one to one and differentiable on H . Define $\lambda(F) \equiv m_p(\mathbf{h}(F \cap H))$.*

Lemma 11.9.3 *λ is a well defined measure on measurable subsets of U and $\lambda \ll m_p$.*

Proof: Since the E_i are disjoint and h is one to one. $\lambda(\cup_i E_i) \equiv m_p(h(\cup_i E_i \cap H)) = \sum_i m_p(h(E_i \cap H)) = \sum_i \lambda(E_i)$. If $m_p(E) = 0$, then $\lambda(E) \equiv m_p(h(E \cap H)) = 0$ because of Lemma 11.8.1. ■

Since $\lambda \ll m_p$, it follows from the Radon Nikodym theorem of Corollary 10.13.14 that there exists $g \in L^1_{loc}(U)$ such that for F a measurable subset of U ,

$$\lambda(F) = m_p(h(F \cap H)) = \int_F g dm_p \quad (11.11)$$

where $g = 0$ off H . To see that this corollary applies, note that both λ and m_p are finite on compact sets and that every open set is a countable union of compact sets.

Now let F be a Borel set so that $h^{-1}(F) \cap H$ is measurable and plays the role of F in the above. Then

$$\begin{aligned} \lambda(h^{-1}(F)) &\equiv m_p(h(h^{-1}(F) \cap H)) \\ &= \int_U \mathcal{X}_{h^{-1}(F) \cap H}(x) g(x) dm_p(x) = \int_H \mathcal{X}_F(h(x)) g(x) dm_p(x) \end{aligned}$$

Thus also for s a Borel measurable nonnegative simple function,

$$\int_{h(H)} s(y) dm_p(y) = \int_H s(h(x)) g(x) dm_p(x)$$

Using a sequence of nonnegative simple functions to approximate a nonnegative Borel measurable f , we obtain from the monotone convergence theorem that

$$\int_{h(H)} f(y) dm_p(y) = \int_H f(h(x)) g(x) dm_p(x)$$

If f is only Lebesgue measurable, then there are nonnegative Borel measurable functions k, l such that $k(y) \leq f(y) \leq l(y)$ with equality holding off a set of m_p measure zero. Then $k(h(x)) g(x) \leq f(h(x)) g(x) \leq l(h(x)) g(x)$ and the two on the ends are Lebesgue measurable which forces the function in the center to also be Lebesgue measurable by completeness of Lebesgue measure because

$$\begin{aligned} \int_H l(h(x)) g(x) - k(h(x)) g(x) dm_p &= \int_{h(H)} l(y) dm_p - \int_{h(H)} k(y) dm_p \\ &= \int_{h(H)} f(y) dm_p - \int_{h(H)} f(y) dm_p = 0 \end{aligned}$$

Thus $l(h(x)) g(x) - k(h(x)) g(x) = 0$ a.e. Then for f nonnegative and Lebesgue measurable,

$$\int_H f(h(x)) g(x) dm_p = \int_{h(H)} f(y) dm_p.$$

This shows the following lemma.

Lemma 11.9.4 *Let $h : U \rightarrow h(U)$ be continuous, U open, and let $H \subseteq U$ be measurable and h is one to one and differentiable on H . Then there exists nonnegative measurable $g \in L^1_{loc}$ such that whenever f is nonnegative and Lebesgue measurable,*

$$\int_{h(H)} f(y) dm_p = \int_H f(h(x)) g(x) dm_p$$

where all necessary measurability is obtained.

It remains to identify g .

Lemma 11.9.5 *For a.e. x , satisfying $|\det D\mathbf{h}(x)| > 0$, and r small enough,*

$$\begin{aligned} D\mathbf{h}(x)B(\mathbf{0}, (1-\varepsilon)r) &\subseteq \mathbf{h}(B(x, r)) \subseteq \mathbf{h}(\overline{B(x, r)}) \subseteq D\mathbf{h}(x)\overline{B(\mathbf{0}, (1+\varepsilon)r)}, \\ \frac{m_p(\mathbf{h}(B(x, r)))}{m_p(B(x, r))} &\in [|\det D\mathbf{h}(x)|(1-\varepsilon)^p, |\det D\mathbf{h}(x)|(1+\varepsilon)^p] \\ \lim_{r \rightarrow 0} \frac{m_p(\mathbf{h}(B(x, r)))}{m_p(B(x, r))} &= |\det D\mathbf{h}(x)| \end{aligned}$$

Proof: For r small enough,

$$\begin{aligned} \mathbf{h}(B(x, r)) &\subseteq \mathbf{h}(x) + D\mathbf{h}(x)B(\mathbf{0}, r) + D\mathbf{h}(x)D\mathbf{h}(x)^{-1}B(\mathbf{0}, \varepsilon r) \\ &\subseteq \mathbf{h}(x) + D\mathbf{h}(x)B(\mathbf{0}, r) + D\mathbf{h}(x)B(\mathbf{0}, \varepsilon r) \\ &\subseteq \mathbf{h}(x) + D\mathbf{h}(x)(B(\mathbf{0}, (1+\varepsilon)r)) \end{aligned}$$

and so $m_p(\mathbf{h}(B(x, r))) \leq |\det(D\mathbf{h}(x))| m_p(B(\mathbf{0}, (1+\varepsilon)r))$. Also,

$$\mathbf{h}(x + v) = \mathbf{h}(x) + D\mathbf{h}(x)v + D\mathbf{h}(x)D\mathbf{h}(x)^{-1}\mathbf{o}(v)$$

and so $\|D\mathbf{h}(x)^{-1}(\mathbf{h}(x + v) - \mathbf{h}(x)) - v\| = \|D\mathbf{h}(x)^{-1}\mathbf{o}(v)\| = \|\mathbf{o}(v)\|$. Thus if r is chosen sufficiently small, it follows that for $v \in B(\mathbf{0}, r)$

$$\|D\mathbf{h}(x)^{-1}(\mathbf{h}(x + v) - \mathbf{h}(x)) - v\| < \varepsilon r$$

and so, from Lemma 8.10.1, $B(\mathbf{0}, (1-\varepsilon)r) \subseteq D\mathbf{h}(x)^{-1}(\mathbf{h}(x + \overline{B(\mathbf{0}, r)}) - \mathbf{h}(x))$.

$$\mathbf{h}(\overline{B(x, r)}) = \mathbf{h}(x + \overline{B(\mathbf{0}, r)}) - \mathbf{h}(x) \supseteq D\mathbf{h}(x)B(\mathbf{0}, (1-\varepsilon)r)$$

Therefore, since $m_p(B(x, r)) = m_p(\overline{B(x, r)})$,

$$|\det(D\mathbf{h}(x))| m_p(B(\mathbf{0}, (1-\varepsilon)r)) = |\det(D\mathbf{h}(x))| (1-\varepsilon)^p r^p \alpha_p \leq m_p(\mathbf{h}(B(x, r)))$$

so for r small enough,

$$\frac{m_p(\mathbf{h}(B(x, r)))}{m_p(B(\mathbf{0}, (1+\varepsilon)r))} \leq |\det(D\mathbf{h}(x))| \leq \frac{m_p(\mathbf{h}(B(x, r)))}{m_p(B(\mathbf{0}, (1-\varepsilon)r))}$$

The claim follows from this since $\varepsilon > 0$ is arbitrary. ■

Lemma 11.9.6 *For a.e. x with $|\det D\mathbf{h}(x)| > 0$, $\lim_{r \rightarrow 0} \frac{m_p(\mathbf{h}(B(x, r) \cap H))}{m_p(\mathbf{h}(B(x, r)))} = \frac{g(x)}{|\det D\mathbf{h}(x)|}$.*

Proof: Using the result of Lemma 11.9.5, for a.e. x satisfying $|\det D\mathbf{h}(x)| > 0$, if r small enough, then

$$m_p(\mathbf{h}(B(x, r))) \in [|\det D\mathbf{h}(x)| m_p(B(x, r))(1-\varepsilon)^p, |\det D\mathbf{h}(x)| m_p(B(x, r))(1+\varepsilon)^p]$$

Therefore, for $Q_r \equiv \frac{m_p(\mathbf{h}(B(\mathbf{x}, r) \cap H))}{m_p(\mathbf{h}(B(\mathbf{x}, r)))} \geq \frac{1}{|\det D\mathbf{h}(\mathbf{x})| m_p(B(\mathbf{x}, r)) (1+\varepsilon)^p} \int_{B(\mathbf{x}, r)} g dm_p$ so

$$\begin{aligned} \frac{1}{m_p(B(\mathbf{0}, r)) (1+\varepsilon)^p} \int_{B(\mathbf{x}, r)} \frac{g}{|\det D\mathbf{h}(\mathbf{x})|} dm_p &\leq Q_r \\ &\leq \frac{1}{m_p(B(\mathbf{0}, r)) (1-\varepsilon)^p} \int_{B(\mathbf{x}, r)} \frac{g}{|\det D\mathbf{h}(\mathbf{x})|} dm_p \end{aligned}$$

and so for Lebesgue points of g , a.e. \mathbf{x} with $|\det D\mathbf{h}(\mathbf{x})| \neq 0$,

$$\frac{1}{(1+\varepsilon)^p} \leq \frac{g(\mathbf{x})}{|\det D\mathbf{h}(\mathbf{x})|} \leq \frac{1}{(1-\varepsilon)^p}$$

Then for such \mathbf{x} , $\frac{1}{(1+\varepsilon)^p} \frac{g}{|\det D\mathbf{h}(\mathbf{x})|} \leq \liminf_{r \rightarrow 0} Q_r \leq \limsup_{r \rightarrow 0} Q_r \leq \frac{1}{(1-\varepsilon)^p} \frac{g}{|\det D\mathbf{h}(\mathbf{x})|}$ so, since ε is arbitrary, $\lim_{r \rightarrow 0} Q_r = \frac{g(\mathbf{x})}{|\det D\mathbf{h}(\mathbf{x})|}$. ■

Lemma 11.9.7 For a.e. $\mathbf{x} \in H$, $g(\mathbf{x}) = |\det D\mathbf{h}(\mathbf{x})|$.

Proof: First consider \mathbf{x} such that $|\det(D\mathbf{h}(\mathbf{x}))| \neq 0$. Then by Lemmas 11.9.5 and 11.9.6

$$\begin{aligned} \lim_{r \rightarrow 0} \frac{m_p(\mathbf{h}(B(\mathbf{x}, r) \cap H))}{m_p(B(\mathbf{x}, r))} &= \lim_{r \rightarrow 0} \frac{m_p(\mathbf{h}(B(\mathbf{x}, r) \cap H))}{m_p(\mathbf{h}(B(\mathbf{x}, r)))} \frac{m_p(\mathbf{h}(B(\mathbf{x}, r)))}{m_p(B(\mathbf{x}, r))} \\ &= \frac{g(\mathbf{x})}{|\det D\mathbf{h}(\mathbf{x})|} |\det D\mathbf{h}(\mathbf{x})| = g(\mathbf{x}) \end{aligned}$$

for a.e. \mathbf{x} where $|\det(D\mathbf{h}(\mathbf{x}))| \neq 0$.

If $|\det D\mathbf{h}(\mathbf{x})| = 0$ then for r small enough,

$$\begin{aligned} \frac{1}{m_p(B(\mathbf{x}, r))} \int_{B(\mathbf{x}, r)} g dm_p &= \frac{m_p(\mathbf{h}(B(\mathbf{x}, r) \cap H))}{m_p(B(\mathbf{x}, r))} \\ &\leq \frac{m_p(\mathbf{h}(\mathbf{x}) + D\mathbf{h}(\mathbf{x})B(\mathbf{0}, r) + B(\mathbf{0}, \varepsilon r))}{m_p(B(\mathbf{x}, r))} = \frac{m_p(D\mathbf{h}(\mathbf{x})B(\mathbf{0}, r) + B(\mathbf{0}, \varepsilon r))}{m_p(B(\mathbf{x}, r))} \end{aligned}$$

Now $D\mathbf{h}(\mathbf{x})B(\mathbf{0}, r) + B(\mathbf{0}, \varepsilon r)$ has finite diameter and lies in a $p-1$ dimensional subset. Therefore, from Theorem 11.7.4 on linear mappings, there is an orthogonal matrix Q preserving all distances such that

$$|\det Q| m_p(D\mathbf{h}(\mathbf{x})B(\mathbf{0}, r) + B(\mathbf{0}, \varepsilon r)) = m_p(QD\mathbf{h}(\mathbf{x})B(\mathbf{0}, r) + B(\mathbf{0}, \varepsilon r))$$

where $QD\mathbf{h}(\mathbf{x})B(\mathbf{0}, r)$ lies in a ball in \mathbb{R}^{p-1} of some radius $\hat{r} = \|D\mathbf{h}(\mathbf{x})\| r$. Thus the set on the right side is contained in a cylinder of radius $\hat{r} + \varepsilon r$ and height $2\varepsilon r$ so its measure is no more than $\alpha_{p-1}(\hat{r} + r\varepsilon)^{p-1} 2\varepsilon r$ for $\alpha_{p-1} = m_{p-1}(B(\mathbf{0}, 1))$. Thus,

$$\begin{aligned} \frac{1}{m_p(B(\mathbf{x}, r))} \int_{B(\mathbf{x}, r)} g dm_p &\leq \frac{(\|D\mathbf{h}(\mathbf{x})\| + 1)^p \alpha_{p-1} (r + r\varepsilon)^{p-1} 2\varepsilon r}{\alpha_p r^p} \\ &= 2(\|D\mathbf{h}(\mathbf{x})\| + 1)^p \frac{\alpha_{p-1}}{\alpha_p} (1 + \varepsilon)^{p-1} \varepsilon \end{aligned}$$

Since ε is arbitrary, for every Lebesgue point where $|\det D\mathbf{h}(\mathbf{x})| = 0$, it follows $g = 0 = |\det D\mathbf{h}(\mathbf{x})|$. ■

Here is the change of variables formula which follows from Lemma 11.9.4 now that g has been identified.

Theorem 11.9.8 *Let U be an open set and let $\mathbf{h} : U \rightarrow \mathbf{h}(U)$ be continuous and one to one and differentiable on the measurable $H \subseteq U$. Then if $f \geq 0$ is Lebesgue measurable,*

$$\int_{\mathbf{h}(H)} f(\mathbf{y}) dm_p = \int_H f(\mathbf{h}(x)) |\det(D\mathbf{h}(x))| dm_p$$

11.10 Mappings Not One to One

Let $H \subseteq U$ and \mathbf{h} is differentiable on H . Let $Z \equiv \{x \in H : |\det D\mathbf{h}(x)| = 0\}$. Then it is possible to decompose $H \setminus Z$ into countably many disjoint measurable sets such that \mathbf{h} will be one to one on each of these sets.

Lemma 11.10.1 *Let $\mathbf{h} : H \subseteq U \subseteq \mathbb{R}^p \rightarrow \mathbb{R}^p$ be differentiable on H and let*

$$Z \equiv \{x : \det(D\mathbf{h}(x)) = 0\}.$$

Then there is a countable set $\{F_i\}_i$ of disjoint Borel sets such that \mathbf{h} is one to one on each of these and $[\det(D\mathbf{h}(x)) \neq 0] = H \setminus Z = \cup_i F_i$.

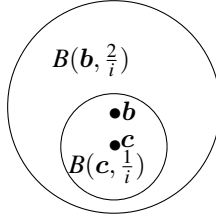
Proof: Let \mathbf{h} be differentiable on the measurable set $H \subseteq U \subseteq \mathbb{R}^p$ for U an open set in \mathbb{R}^p . Let \mathcal{S} be a countable dense subset of the set of invertible matrices and let \mathcal{C} be a countable dense subset of B , a Borel subset of the points x of U where $\det(D\mathbf{h}(x)) \neq 0$. To get \mathcal{S} one could simply consider all matrices of which have a rational number in each entry. This would be dense in $\mathcal{L}(\mathbb{R}^p, \mathbb{R}^p)$ which is therefore a separable metric space, which therefore has a countable basis of open balls. Then \mathcal{S} being a subset must also be separable. (Corollary 3.4.3) I will decompose B into a disjoint union of Borel sets on which \mathbf{h} is one to one. This will be done by establishing 11.15 given below where T is an invertible transformation. For $T \in \mathcal{S}$, $c \in \mathcal{C}$, $i \in \mathbb{N}$, define $E(T, c, i)$ to be those $b \in B(c, \frac{1}{i})$ such that for all $a \in B(b, \frac{2}{i})$,

$$|\mathbf{h}(a) - \mathbf{h}(b) - D\mathbf{h}(b)(a - b)| < \varepsilon |T(a - b)| \quad (11.12)$$

and also $D\mathbf{h}(b)$ is close enough to T that the following hold.

$$\inf_{v \neq 0} \frac{|D\mathbf{h}(b)v|}{|Tv|} > 1 - \varepsilon, \quad \sup_{v \neq 0} \frac{|D\mathbf{h}(b)v|}{|Tv|} < 1 + \varepsilon \quad (11.13)$$

where $\varepsilon < 1/4$. These are Borel sets because of continuity of \mathbf{h} so that the derivative $D\mathbf{h}$ is also Borel measurable. Indeed each entry of the matrix of $D\mathbf{h}$ is a limit of difference quotients which are continuous.



Note that it is not clear whether $c \in E(T, c, i)$ because of the above two requirements 11.13. What is going on here is that we are looking for b such that $D\mathbf{h}(b)$ is sufficiently close to one of those T which also are in a piece of B . Thus we start with one of those T

and one of those points c and look for all b , if any, which do the right things. There are countably many of these pieces of B being denoted as $E(T, c, i)$.

The union of these $E(T, c, i)$ is all of B because if $b \in B$,

$$|h(a) - h(b) - Dh(b)(a - b)| < \varepsilon |Dh(b)(a - b)| \quad (11.14)$$

whenever $a \in B(b, \frac{2}{i})$ provided i is sufficiently large. Thus also, by Lemma 5.3.1, there is $T \in \mathcal{S}$ such that the above holds for $Dh(b)$ replaced with T and $a \in B(b, \frac{2}{i})$ and also 11.13. Thus $b \in E(T, c, i)$, so indeed the union of these sets is B .

Now let $a, b \in E(T, c, i)$. Since $a, b \in E(T, c, i)$, a, b are within $1/i$ of c and so a is within $2/i$ of b and so 11.12 holds because of the definition of $E(T, c, i)$. Therefore, from 11.12 and the inequalities which follow, 11.13,

$$(1 - 3\varepsilon) |T(a - b)| \leq |h(a) - h(b)| \leq (1 + 3\varepsilon) |T(a - b)| \quad (11.15)$$

Indeed from 11.14 and 11.13

$$|h(a) - h(b)| \leq (1 + \varepsilon) |Dh(b)(a - b)| \leq (1 + \varepsilon)^2 |T(a - b)| \leq (1 + 3\varepsilon) |T(a - b)|$$

The bottom inequality is similar. Thus h is one to one on $E(T, c, i)$. Now enumerate these Borel sets $\{E_i\}_{i=1}^\infty$. Let $F_1 = E_1$ and if F_1, \dots, F_m have been chosen, let $F_{m+1} \equiv E_{m+1} \setminus (\cup_{i=1}^m F_i)$. ■

Theorem 11.10.2 *Let U be an open set and let $h : U \rightarrow h(U)$ be continuous and differentiable on the measurable $H \subseteq U$ such that $h(U \setminus H)$ has measure zero. Then if $f \geq 0$ is Lebesgue measurable,*

$$\int_{h(H)} \#(y) f(y) dm_p = \int_H f(h(x)) |\det(Dh(x))| dm_p$$

where $\#(y)$ is the number of elements of $h^{-1}(y)$ in U .

Proof: Let $\{F_i\}$ be the Borel sets of Lemma 11.10.1 whose union equals $H \setminus Z$ where Z is the set where $Dh(x)$ exists but is not invertible and h one to one on each F_i . Thus for f Lebesgue measurable, $\int_{h(H \setminus Z)} \mathcal{R}_{h(F_i)} f(y) dm_p = \int_{F_i} f(h(x)) |\det(Dh(x))| dm_p$. Let $n(y) \equiv \sum_i \mathcal{R}_{h(F_i)}(y)$. Then, adding these yields,

$$\int_{h(H \setminus Z)} n(y) f(y) dm_p = \int_{H \setminus Z} f(h(x)) |\det(Dh(x))| dm_p$$

Now $\#(y) = n(y)$ except for $h(U \setminus H) \cup h(Z)$ which is a set of m_p measure zero by assumption and by Sard's Lemma, Lemma 11.8.3 for $h(Z)$. Therefore,

$$\begin{aligned} \int_{h(H)} \#(y) f(y) dm_p &= \int_{h(H \setminus Z)} \#(y) f(y) dm_p = \int_{H \setminus Z} f(h(x)) |\det(Dh(x))| dm_p \\ &= \int_H f(h(x)) |\det(Dh(x))| dm_p \quad \blacksquare \end{aligned}$$

h is one to one when $\#(y) = 1$ and in this case we get the usual change of variables formula.

11.11 Spherical Coordinates

As usual, S^{p-1} is the unit sphere, the boundary of the unit ball $B(\mathbf{0}, 1)$ in \mathbb{R}^p . It is a metric space with respect to the usual notion of distance which it inherits from being a part of \mathbb{R}^p . Then $(0, \infty) \times S^{p-1}$ is also a metric space with the metric

$$d((\rho, \omega), (\hat{\rho}, \hat{\omega})) \equiv \max\{|\rho - \hat{\rho}|, |\omega - \hat{\omega}|\}$$

Indeed, this kind of thing delivers a metric for an arbitrary finite product of metric spaces. See Problem 6 on Page 94.

Definition 11.11.1 Define $\lambda : \mathbb{R}^p \setminus \{\mathbf{0}\} \rightarrow (0, \infty) \times S^{p-1}$ as $\lambda(x) \equiv \left(|x|, \frac{x}{|x|}\right)$

Then with this definition, the following is true.

Lemma 11.11.2 Let λ be as defined above. Then λ is one to one, onto, and continuous with continuous inverse.

Proof: First of all, it is obviously onto. Indeed, if $(\rho, \omega) \in (0, \infty) \times S^{p-1}$, consider $x \equiv \rho\omega$. Why is this one to one? If $x \neq \hat{x}$, then there are two cases. It might be that $|x| \neq |\hat{x}|$ and in this case, clearly $\lambda(x) \neq \lambda(\hat{x})$. The other case is that $|x| = |\hat{x}| = \rho$ but these two vectors x, \hat{x} are not equal. In this case, $\frac{x}{|x|} - \frac{\hat{x}}{|\hat{x}|} = \frac{1}{\rho}(x - \hat{x}) \neq \mathbf{0}$. Thus λ is one to one.

Is λ continuous? Suppose $x_n \rightarrow x \neq \mathbf{0}$. Does $\lambda(x_n) \rightarrow \lambda(x)$? First of all, the triangle inequality shows that $|x_n| \rightarrow |x|$. It only remains to verify $\frac{x_n}{|x_n|} \rightarrow \frac{x}{|x|}$. This is clearly the case because

$$\left| \frac{x_n}{|x_n|} - \frac{x}{|x|} \right| = \left| \frac{|x|x_n - |x_n|x}{|x_n||x|} \right| \rightarrow \left| \frac{|x|x - |x|x}{|x||x|} \right| = 0$$

Is λ^{-1} also continuous? One could show this directly or observe that λ^{-1} is automatically continuous on $[\frac{1}{n}, n] \times S^{p-1}$ because this is a compact set. Indeed, $[\frac{1}{n}, n] \times S^{p-1} = \lambda(\{x \in \mathbb{R}^p : \frac{1}{n} \leq |x| \leq n\})$. If $\lambda(x_n) \rightarrow \lambda(x)$, does it follow that $x_n \rightarrow x$? If not, there exists a subsequence, still denoted as x_n such that $x_n \rightarrow y \neq x$. But then, by continuity of λ , $\lambda(x_n) \rightarrow \lambda(y)$ and so $\lambda(y) = \lambda(x)$ which does not occur because λ is one to one.

It follows, since $(0, \infty) \times S^{p-1} = \cup_n [\frac{1}{n}, n] \times S^{p-1}$, that λ^{-1} is continuous. ■

Thus the open sets for $(0, \infty) \times S^{p-1}$ are all of the form $\lambda(U)$ where U is open in $\mathbb{R}^p \setminus \{\mathbf{0}\}$. Also, the open sets of $\mathbb{R}^p \setminus \{\mathbf{0}\}$ are of the form $\lambda^{-1}(V)$ where V is an open set of $(0, \infty) \times S^{p-1}$. One can replace the word “open” with the word “Borel” in the previous observation.

Motivated by familiar formulas for the area of a sphere and the circumference of a circle, here is a definition of a surface measure defined on the Borel sets of S^{p-1} .

Definition 11.11.3 Let E be a Borel set on S^{p-1} . Then

$$\lambda^{-1}((0, 1] \times E) \equiv \{\rho\omega : \rho \in (0, 1], \omega \in E\}$$

will be a part of the unit ball formed from the cone starting at $\mathbf{0}$ and extending to the points of E , leaving out $\mathbf{0}$. Since $(0, 1] \times E$ is a Borel set in $(0, \infty) \times S^{p-1}$ thanks to Problem 4 on Page 259, this cone just described is a Borel set in \mathbb{R}^p . Then

$$\sigma(E) \equiv pm_p\left(\lambda^{-1}((0, 1] \times E)\right)$$

This is obviously a measure on the Borel sets of S^{p-1} .

Is this even a good idea? Note $m_p(\lambda^{-1}(\{r\} \times E)) = 0$ because $\lambda^{-1}(\{r\} \times E)$ is just a part of the sphere of radius r which has m_p measure zero. The reason this is so is as follows. Letting $\alpha_p \equiv m_p(B(\mathbf{0}, 1))$, the sphere of radius r is contained in $B(\mathbf{0}, r + \varepsilon) \setminus B(\mathbf{0}, r - \varepsilon)$ and so the sphere has m_p measure no more than $\alpha_p((r + \varepsilon)^p - (r - \varepsilon)^p)$ for every $\varepsilon > 0$.

Lemma 11.11.4 *Let G be a Borel set in $(0, \infty) \times S^{p-1}$. Then*

$$m_p(\lambda^{-1}(G)) = \int_0^\infty \int_{S^{p-1}} \mathcal{X}_G(\rho, \omega) \rho^{p-1} d\sigma d\rho \quad (11.16)$$

and the iterated integrals make sense.

Proof: Let $\mathcal{K} \equiv \{I \times E : I \text{ is an interval in } (0, \infty) \text{ and } E \text{ is Borel in } S^{p-1}\}$. This is a π system and each set of \mathcal{K} is a Borel set. Then if I is one of these intervals, having end points $a < b$,

$$\begin{aligned} \int_0^\infty \int_{S^{p-1}} \mathcal{X}_{I \times E}(\rho, \omega) \rho^{p-1} d\sigma d\rho &= \int_a^b \rho^{p-1} \int_E d\sigma d\rho = \sigma(E) \left(\frac{b^p}{p} - \frac{a^p}{p} \right) \\ &= p m_p(\lambda^{-1}((0, 1] \times E)) \left(\frac{b^p}{p} - \frac{a^p}{p} \right) = m_p(\lambda^{-1}((0, 1] \times E)) (b^p - a^p) \\ &= m_p(\lambda^{-1}((a, b) \times E)) = m_p(\lambda^{-1}(I \times E)) \end{aligned}$$

Let \mathcal{G} denote those Borel sets G in $(0, \infty) \times S^{p-1}$ for which, $G_n \equiv G \cap (0, n) \times S^{p-1}$,

$$m_p(\lambda^{-1}(G_n)) = \int_0^\infty \int_{S^{p-1}} \mathcal{X}_{G_n}(\rho, \omega) \rho^{p-1} d\sigma d\rho$$

the iterated integrals making sense. It is routine to verify that \mathcal{G} is closed with respect to complements and countable disjoint unions. It was also shown above that it contains \mathcal{K} . By Dynkin's lemma, Lemma 9.3.2, \mathcal{G} equals the Borel sets in $(0, \infty) \times S^{p-1}$. Now use the monotone convergence theorem. ■

Theorem 11.11.5 *Let f be a Borel measurable nonnegative function. Then*

$$\int f dm_p = \int_0^\infty \int_{S^{p-1}} f(\rho \omega) \rho^{p-1} d\sigma d\rho \quad (11.17)$$

Proof: From the above lemma, if F is an arbitrary Borel set, it has the same measure as $F \cap (\mathbb{R}^p \setminus \{\mathbf{0}\})$ so there is no loss of generality in assuming $\mathbf{0} \notin F$.

$$\begin{aligned} \int_{\mathbb{R}^p} \mathcal{X}_F dm_p &= m_p(F) = m_p(\lambda^{-1}(\lambda(F))) = \int_0^\infty \int_{S^{p-1}} \mathcal{X}_{\lambda(F)}(\rho, \omega) \rho^{p-1} d\sigma d\rho \\ &= \int_0^\infty \int_{S^{p-1}} \mathcal{X}_F(\lambda^{-1}(\rho, \omega)) \rho^{p-1} d\sigma d\rho = \int_0^\infty \int_{S^{p-1}} \mathcal{X}_F(\rho \omega) \rho^{p-1} d\sigma d\rho \end{aligned}$$

Now if f is nonnegative and Borel measurable, one can approximate using Borel simple functions increasing pointwise to f and use the monotone convergence theorem to obtain 11.17. ■

Note that by Theorem 10.14.9, you can interchange the order of integration in 11.16 if desired.

Example 11.11.6 For what values of s is the integral $\int_{B(\mathbf{0}, R)} (1 + |\mathbf{x}|^2)^s dy$ bounded independent of R ? Here $B(\mathbf{0}, R)$ is the ball, $\{\mathbf{x} \in \mathbb{R}^p : |\mathbf{x}| \leq R\}$.

I think you can see immediately that s must be negative but exactly how negative? It turns out it depends on p and using polar coordinates, you can find just exactly what is needed. From the polar coordinates formula above,

$$\int_{B(\mathbf{0}, R)} (1 + |\mathbf{x}|^2)^s dy = \int_0^R \int_{S^{p-1}} (1 + \rho^2)^s \rho^{p-1} d\sigma d\rho = C_p \int_0^R (1 + \rho^2)^s \rho^{p-1} d\rho$$

Now the very hard problem has been reduced to considering an easy one variable problem of finding when $\int_0^R \rho^{p-1} (1 + \rho^2)^s d\rho$ is bounded independent of R . You need $2s + (p-1) < -1$ so you need $s < -p/2$.

11.12 Symmetric Derivative for Radon Measures

Here we have two Radon measures μ, λ defined on a σ algebra of sets \mathcal{F} which are subsets of an open subset U of \mathbb{R}^p , possibly all of \mathbb{R}^p . They are complete and Borel and inner and outer regular, and finite on compact sets. Thus both of these measures are σ finite.

In this section is the symmetric derivative $D_\mu(\lambda)$. In what follows, $B(\mathbf{x}, r)$ will denote a **closed** ball with center \mathbf{x} and radius r . Also, let λ and μ be Radon measures and as above, Z will denote a μ measure zero set off of which $\mu(B(\mathbf{x}, r)) > 0$ for all $r > 0$. Generalizing the notion of \limsup and \liminf ,

$$\limsup_{r \rightarrow 0} f(r) \equiv \lim_{r \rightarrow 0} (\sup \{f(t) : t < r\}), \quad \liminf_{r \rightarrow 0} f(r) \equiv \lim_{r \rightarrow 0} (\inf \{f(t) : t < r\})$$

Then directly from this definition, the $\lim_{r \rightarrow 0}$ exists if and only if these two are equal.

Definition 11.12.1 For $\mathbf{x} \notin Z$, define the upper and lower symmetric derivatives as

$$\overline{D}_\mu \lambda(\mathbf{x}) \equiv \limsup_{r \rightarrow 0} \frac{\lambda(B(\mathbf{x}, r))}{\mu(B(\mathbf{x}, r))}, \quad \underline{D}_\mu \lambda(\mathbf{x}) \equiv \liminf_{r \rightarrow 0} \frac{\lambda(B(\mathbf{x}, r))}{\mu(B(\mathbf{x}, r))}.$$

respectively. Also define $D_\mu \lambda(\mathbf{x}) \equiv \overline{D}_\mu \lambda(\mathbf{x}) = \underline{D}_\mu \lambda(\mathbf{x})$ in the case when both the upper and lower derivatives are equal. Recall that $Z \equiv \{\mathbf{x} : \mu(B(\mathbf{x}, r)) = 0 \text{ for some } r > 0\}$ and that this set has measure zero.

Lemma 11.12.2 Let λ and μ be Radon measures on \mathcal{F}_λ and \mathcal{F}_μ respectively and let $a, b > 0$. If A is a subset of $\{\mathbf{x} \notin Z : \overline{D}_\mu \lambda(\mathbf{x}) \geq b\}$ then $\overline{\lambda}(A) \geq b\overline{\mu}(A)$ and if A is a subset of $\{\mathbf{x} \notin Z : \underline{D}_\mu \lambda(\mathbf{x}) \leq a\}$, then $\overline{\lambda}(A) \leq a\overline{\mu}(A)$.

Proof: Let $\overline{\lambda}$ be the outer measure determined by λ , similar for $\overline{\mu}$ and μ . Suppose first that A is a subset of $\{\mathbf{x} \notin Z : \overline{D}_\mu \lambda(\mathbf{x}) \geq b\}$ so $\mu(B(\mathbf{x}, r)) > 0$ for all $r > 0$ and $\overline{\lambda}(A) < \infty$. Let small $\varepsilon > 0$, and let V be a bounded open set with $V \supseteq A$ and $\lambda(V) - \varepsilon < \overline{\lambda}(A)$. Then for each $\mathbf{x} \in A$, $\frac{\lambda(B(\mathbf{x}, r))}{\mu(B(\mathbf{x}, r))} > b - \varepsilon$, $B(\mathbf{x}, r) \subseteq V$, for infinitely many values of r which are arbitrarily small. Thus the collection of such closed balls constitutes a Vitali cover for A . By Corollary 9.12.3 there is a disjoint sequence of these closed balls $\{B_i\}$ such that $\overline{\mu}(A \setminus \bigcup_{i=1}^\infty B_i) = 0$,

$$\overline{\mu}(A) \leq \overline{\mu}(A \setminus \bigcup_{i=1}^\infty B_i) + \overline{\mu}(\bigcup_{i=1}^\infty B_i \cap A) \leq \sum_{i=1}^\infty \overline{\mu}(B_i \cap A) \quad (11.18)$$

$$\leq \sum_{i=1}^{\infty} \mu(B_i) < \frac{1}{b-\varepsilon} \sum_{i=1}^{\infty} \lambda(B_i) \leq \frac{1}{b-\varepsilon} \lambda(V) < \frac{\varepsilon + \bar{\lambda}(A)}{b-\varepsilon}$$

Since ε is arbitrary, this shows $\bar{\mu}(A)b \leq \bar{\lambda}(A)$. In case $\bar{\lambda}(A) = \infty$, there is nothing to show.

Now suppose A is a subset of $\{x \notin Z : \underline{D}_\mu \lambda(x) \leq a\}$. Suppose first that A is bounded so $\bar{\mu}(A) < \infty$ and let V be a bounded open set containing A with $\mu(V) - \varepsilon < \bar{\mu}(A)$. Then for each $x \in A$, there are arbitrarily small r such that $\frac{\lambda(B(x,r))}{\mu(B(x,r))} < a + \varepsilon$, $B(x,r) \subseteq V$. Therefore, by Corollary 9.12.3 again, there exists a disjoint sequence of these balls, $\{B_i\}$ satisfying this time,

$$\begin{aligned} \bar{\lambda}(A \setminus \bigcup_{i=1}^{\infty} B_i) &= 0 \text{ and } \bar{\lambda}(A) \leq \sum_{i=1}^{\infty} \bar{\lambda}(A \cap B_i) \leq \sum_{i=1}^{\infty} \lambda(B_i) \\ &\leq \sum_{i=1}^{\infty} (a + \varepsilon) \mu(B_i) \leq (a + \varepsilon) \mu(V) \leq (a + \varepsilon)(\varepsilon + \bar{\mu}(A)) \end{aligned}$$

Since $\varepsilon > 0$ is arbitrary, $\bar{\lambda}(A) \leq a\bar{\mu}(A)$. This proves the lemma in case A is bounded. If not, replace A with $A \cap B(0, R)$, get the result for this and let $R \rightarrow \infty$. If $\bar{\mu}(A) = \infty$, there is nothing to show. ■

Theorem 11.12.3 *Let λ, μ be Radon measures on \mathcal{F}_λ and \mathcal{F}_μ respectively. There exists a set of measure zero N containing Z such that for $x \notin N$, $D_\mu \lambda(x)$ exists and also $\mathcal{X}_{N^c}(\cdot) D_\mu \lambda(\cdot)$ is a μ measurable function. Furthermore, $D_\mu \lambda(x) < \infty$ μ a.e.*

Proof: First I show $D_\mu \lambda(x)$ exists a.e. Let $0 < a < b < \infty$ and let A be any bounded subset of $N(a, b) \equiv \{x \notin Z : \bar{D}_\mu \lambda(x) > b > a > \underline{D}_\mu \lambda(x)\}$. By Lemma 11.12.2, $a\bar{\mu}(A) \geq \bar{\lambda}(A) \geq b\bar{\mu}(A)$ and so $\bar{\mu}(A) = 0$ and so A is μ measurable. It follows $\mu(N(a, b)) = 0$ because $\mu(N(a, b)) \leq \sum_{m=1}^{\infty} \mu(N(a, b) \cap B(0, m)) = 0$. Now

$$\{x \notin Z : \bar{D}_\mu \lambda(x) > \underline{D}_\mu \lambda(x)\} \subseteq \bigcup \{N(a, b) : 0 < a < b, \text{ and } a, b \in \mathbb{Q}\}$$

and the latter set is a countable union of sets of measure 0, so off a set of measure 0, N for which $N \supseteq Z$, one has $\bar{D}_\mu \lambda(x) = \underline{D}_\mu \lambda(x)$.

We can assume also that N is a Borel set from regularity considerations. See Theorem 9.11.2 for example. It remains to verify $\mathcal{X}_{N^c}(\cdot) D_\mu \lambda(\cdot)$ is finite a.e. and is μ measurable. Let $I = \{x : D_\mu \lambda(x) = \infty\}$. Then by Lemma 11.12.2 $\bar{\lambda}(I \cap B(0, m)) \geq a\bar{\mu}(I \cap B(0, m))$ for all $a > 0$, and since λ is finite on bounded sets, the above implies $\bar{\mu}(I \cap B(0, m)) = 0$ for each m which implies that I is μ measurable and has μ measure zero since $I = \bigcup_{m=1}^{\infty} I \cap B(0, m)$.

Now the issue is measurability. Let λ be an arbitrary Radon measure. I need show that $x \rightarrow \lambda(B(x, r))$ is measurable. Here is where it is convenient to have the balls be closed balls. If V is an open set containing $B(x, r)$, then for y close enough to x , $B(y, r) \subseteq V$ also and so, $\limsup_{y \rightarrow x} \lambda(B(y, r)) \leq \lambda(V)$. However, since V is arbitrary and λ is outer regular, or observing that $B(x, r)$ the closed ball is the intersection of nested open sets, it follows that $\limsup_{y \rightarrow x} \lambda(B(y, r)) \leq \lambda(B(x, r))$. Thus $x \rightarrow \lambda(B(x, r))$ is upper semicontinuous, similar for $x \rightarrow \mu(B(x, r))$ and so, $x \rightarrow \frac{\lambda(B(x, r))}{\mu(B(x, r))}$ is measurable. Hence

$$\mathcal{X}_{N^c}(x) D_\mu(\lambda)(x) = \lim_{r_i \rightarrow 0} \mathcal{X}_{N^c}(x) \frac{\lambda(B(x, r_i))}{\mu(B(x, r_i))} \text{ is also measurable. } \blacksquare$$

Typically I will write $D_\mu \lambda(x)$ rather than the more precise $\mathcal{X}_{N^c}(x) D_\mu \lambda(x)$ since the values on the set of measure zero N are not important due to the completeness of the measure μ . This is done in the next section.

11.13 Radon Nikodym Theorem, Radon Measures

The Radon Nikodym theorem is an abstract result but this will be a special version. It will give a pointwise description in terms of the symmetric derivative of the Radon Nikodym derivative presented earlier.

Definition 11.13.1 Let λ, μ be two Radon measures defined on \mathcal{F} , a σ algebra of subsets of an open set U . Then $\lambda \ll \mu$ means that whenever $\mu(E) = 0$, it follows that $\lambda(E) = 0$.

Next is a representation theorem for λ in terms of an integral involving $D_\mu \lambda$.

Theorem 11.13.2 Let λ and μ be Radon measures defined on $\mathcal{F}_\lambda, \mathcal{F}_\mu$ respectively, σ algebras of the open set U , then there exists a set of μ measure zero N such that $D_\mu \lambda(x)$ exists off N and if $E \subseteq N^C, E \in \mathcal{F}_\lambda \cap \mathcal{F}_\mu$, then $\lambda(E) = \int_U (D_\mu \lambda) \mathcal{H}_E d\mu$. If $\lambda \ll \mu$ on $\mathcal{F}_\lambda \cap \mathcal{F}_\mu$, then $\lambda(E) = \int_E D_\mu \lambda d\mu$. In any case, $\lambda(E) \geq \int_E D_\mu \lambda d\mu$ so $D_\mu \lambda$ is in $L^1_{loc}(\mathbb{R}^p, \mu)$ because $\lambda(B) < \infty$ for any ball B .

Proof: The proof is based on Lemma 11.12.2. Let $E \subseteq N^C$ where N has μ measure 0 and includes the set Z along with the set where the symmetric derivative does not exist. It can be assumed that N is a G_δ set. Define

$$l_n(x) \equiv \sum_{k=1}^{\infty} a_{k-1}^n \mathcal{H}_{(D_\mu \lambda)^{-1}(I_k^n)}(x), \quad u_n(x) \equiv \sum_{k=1}^{\infty} a_k^n \mathcal{H}_{(D_\mu \lambda)^{-1}(I_k^n)}(x)$$

where $I_k^n \equiv ((k-1)2^{-n}, k2^{-n}] \equiv (a_{k-1}^n, a_k^n]$ for $k, n \in \mathbb{N}$. Thus $u_n(x) \geq D_\mu \lambda(x) > l_n(x)$ and $u_n(x) - l_n(x) = 2^{-n}$. Also, $l_n(x)$ increases to $D_\mu \lambda(x)$. Letting

$$E_k^n \equiv [x \in E : D_\mu \lambda(x) \in I_k^n],$$

and assuming $\mu(E) < \infty, \int_E D_\mu \lambda d\mu \in [\int_E l_n d\mu, \int_E u_n d\mu]$

$$= \left[\sum_{k=1}^{\infty} a_{k-1}^n \mu(E_k^n), \sum_{k=1}^{\infty} a_k^n \mu(E_k^n) \right] \subseteq \left[\int_E l_n d\mu, \int_E l_n d\mu + 2^{-n} \mu(E) \right] \quad (11.19)$$

From Lemma 11.12.2, $\mu(E_k^n) a_k^n \geq \lambda(E_k^n) \geq a_{k-1}^n \mu(E_k^n)$ and so the interval in 11.19 contains $\sum_{k=1}^{\infty} \lambda(E_k^n)$. This equals $\lambda(E)$ because of Lemma 11.12.2 which implies

$$\begin{aligned} \lambda(E \cap \{x \in N^C : D_\mu \lambda(x) = 0\}) &\leq a\mu(E \cap \{x \in N^C : D_\mu \lambda(x) = 0\}) \\ &\leq a\mu(E), \quad \mu(E) < \infty \end{aligned}$$

and since this is true for every positive a , it follows that

$$\lambda(E \cap \{x \in N^C : D_\mu \lambda(x) = 0\}) = 0$$

so the sum $\sum_{k=1}^{\infty} \lambda(E_k^n) = \lambda(E)$. Then, from the monotone convergence theorem in 11.19, one can pass to a limit and find that $\int_E D_\mu \lambda d\mu = \lambda(E)$.

Now if E is an arbitrary set in N^C , maybe not bounded, the above shows

$$\lambda(E \cap B(\mathbf{0}, n)) = \int_{E \cap B(\mathbf{0}, n)} D_\mu \lambda d\mu$$

Let $n \rightarrow \infty$ and use the monotone convergence theorem. Thus for all $E \subseteq N^C$, $\lambda(E) = \int_E D_\mu \lambda d\mu$. For the last claim, $\int_E D_\mu \lambda d\mu = \int_{E \cap N^C} D_\mu \lambda d\mu = \lambda(E \cap N^C) \leq \lambda(E)$.

In case, $\lambda \ll \mu$, it does not matter that $E \subseteq N^C$ because, since $\mu(N) = 0$, so is $\lambda(N)$ and so

$$\lambda(E) = \lambda(E \cap N^C) = \int_{E \cap N^C} D_\mu \lambda d\mu = \int_E D_\mu \lambda d\mu$$

for any $E \in \mathcal{F}$. ■

What if λ and μ are just two arbitrary Radon measures defined on \mathcal{F} ? What then? It was shown above that $D_\mu \lambda(x)$ exists for μ a.e. x , off a G_δ set N of μ measure 0 which includes Z , the set of x where $\mu(B(x, r)) = 0$ for some $r > 0$. Also, it was shown above that if $E \subseteq N^C$, then $\lambda(E) = \int_E D_\mu \lambda(x) d\mu$. Define for arbitrary $E \in \mathcal{F}$,

$$\lambda_\mu(E) \equiv \lambda(E \cap N^C), \lambda_\perp(E) \equiv \lambda(E \cap N)$$

Then

$$\begin{aligned} \lambda(E) &= \lambda(E \cap N) + \lambda(E \cap N^C) = \lambda_\perp(E) + \lambda_\mu(E) \\ &= \lambda(E \cap N) + \int_{E \cap N^C} D_\mu \lambda(x) d\mu = \lambda(E \cap N) + \int_E D_\mu \lambda(x) d\mu \\ &\equiv \lambda(E \cap N) + \lambda_\mu(E) \equiv \lambda_\perp(E) + \lambda_\mu(E) \end{aligned}$$

This shows the following corollary.

Corollary 11.13.3 *Let μ, λ be two Radon measures. Then there exist two measures, $\lambda_\mu, \lambda_\perp$ such that $\lambda_\mu \ll \mu$, $\lambda = \lambda_\mu + \lambda_\perp$ and a set of μ measure zero N such that $\lambda_\perp(E) = \lambda(E \cap N)$. Also λ_μ is given by the formula $\lambda_\mu(E) \equiv \int_E D_\mu \lambda(x) d\mu$.*

Proof: If $x \in N$, this could happen two ways, either $x \in Z$ or $D_\mu \lambda(x)$ fails to exist. It only remains to verify that λ_μ given above satisfies $\lambda_\mu \ll \mu$. However, this is obvious because if $\mu(E) = 0$, then $\int_E D_\mu \lambda(x) d\mu = 0$. ■

Since $D_\mu \lambda(x) = D_\mu \lambda(x) \mathcal{X}_{N^C}(x)$, it doesn't matter which we use but maybe $D_\mu \lambda(x)$ doesn't exist at some points of N , so although I will use $D_\mu \lambda(x)$, it might be more precise to use $D_\mu \lambda(x) \mathcal{X}_{N^C}(x)$.

This is sometimes called the Lebesgue decomposition.

How does this relate to Corollary 10.13.14? It tells how to find the function f in that Corollary as a symmetric derivative. This is very useful when you want to have an explicit description of the Radon Nikodym derivative.

11.14 Absolutely Continuous Functions

Can you integrate the derivative to get the function as in calculus? The answer is that sometimes you can and when this is the case, the function is called absolutely continuous. This is explained in this section. Recall the following which summarizes Theorems 9.9.1 on Page 257 and 9.7.4 on Page 250. In what follows m will be one dimensional Lebesgue measure. Recall that for F increasing, $F(x+) \equiv \lim_{h \rightarrow 0+} F(x+h)$, $F(x-) \equiv \lim_{h \rightarrow 0+} F(x-h)$.

Theorem 11.14.1 *Let F be an increasing function on \mathbb{R} . Then there is an outer measure μ and a σ algebra \mathcal{F} on which μ is a measure such that \mathcal{F} contains the Borel sets. This measure μ satisfies*

$$\mu([a, b]) = F(b+) - F(a-), \mu((a, b)) = F(b-) - F(a+)$$

$$\mu((a, b]) = F(b+) - F(a+), \quad \mu([a, b)) = F(b-) - F(a-).$$

Furthermore, if E is any set in \mathcal{F}

$$\mu(E) = \sup \{ \mu(K) : K \text{ compact, } K \subseteq E \} \quad (11.20)$$

$$\mu(E) = \inf \{ \mu(V) : V \text{ is an open set } V \supseteq E \} \quad (11.21)$$

Of interest is the symmetric derivative of F , defined as $\lim_{h \rightarrow 0} \frac{F((x+h)+) - F((x-h)-)}{2h}$.
Now from Theorem 11.12.3 it follows that for a.e. x ,

$$\lim_{h \rightarrow 0+} \frac{\mu([x-h, x+h])}{2h} = \lim_{h \rightarrow 0+} \frac{F((x+h)+) - F((x-h)-)}{2h} \equiv D_m \mu(x)$$

exists a.e. x . From Corollary 11.13.3, $\mu = \mu_\perp + \mu_m$ where $\mu_m(E) \equiv \int_E D_m \mu(x) dm$ and there is a Borel set of m measure zero N such that $\mu_\perp(E) = \mu(E \cap N)$

In case $\mu \ll m$ then from Theorem 11.13.2, if E is Borel, $\mu(E) = \int_E D_m \mu(x) dm$. This begs the following question.

What properties on F are equivalent to this measure μ being absolutely continuous with respect to m , $\mu \ll m$? Here is a definition of what it means for a **function** to be absolutely continuous with respect to Lebesgue measure. Thus there are now two things being called “absolutely continuous” functions and measures.

Definition 11.14.2 Let $[a, b]$ be a closed and bounded interval and let $F : [a, b] \rightarrow \mathbb{R}$. Then F is said to be absolutely continuous if for every $\varepsilon > 0$ there exists $\delta > 0$ such that if $\sum_{i=1}^m |y_i - x_i| < \delta$ where the intervals (x_i, y_i) are non-overlapping, then it follows that $\sum_{i=1}^m |F(y_i) - F(x_i)| < \varepsilon$.

It turns out that if, in the definition, you allow arbitrary intervals, non-overlapping or not, then you end up with a Lipschitz function which is obviously absolutely continuous. This is shown later.

The following theorem gives the desired equivalence between absolute continuity of F and $\mu \ll m$.

Theorem 11.14.3 Let F be an increasing function on $[a, b]$ and let μ be the measure of Theorem 11.14.1. Then $\mu \ll m$ if and only if F is absolutely continuous.

Proof: \Rightarrow First suppose that $\mu \ll m$. Then by Theorem 11.13.2, for all Borel sets E ,

$$\mu(E) = \int_E D_m \mu(x) dm$$

In particular, F must be continuous because $D_m \mu$ is in L^1_{loc} . Thus

$$F(y-) - F(x+) = \int_{(x,y)} D_m \mu(x) dm = \int_{(x,y]} D_m \mu(x) dm = F(y+) - F(x+)$$

showing that for arbitrary y , $F(y-) = F(y+)$ so the function F is continuous as claimed. Also $F(b) - F(a) = \int_{[a,b]} D_m \mu(x) dm$ so $D_m \mu$ is in $L^1([a, b], m)$.

If the function F is not absolutely continuous, then there exists $\varepsilon > 0$ and open sets E_n consisting of unions of finitely many non-overlapping open intervals such that if $E_n = \cup_{i=1}^{m_n} (x_i^n, y_i^n)$, then $\sum_{i=1}^{m_n} |y_i^n - x_i^n| = m(E_n) < 2^{-n}$ but

$$\int_{[a,b]} \mathcal{R}_{E_n}(x) D_m \mu(x) dm = \mu(E_n) = \sum_{i=1}^{m_n} \mu(x_i^n, y_i^n) = \sum_{i=1}^{m_n} (F(y_i^n) - F(x_i^n)) \geq \varepsilon \quad (11.22)$$

However, $\mathcal{R}_{E_n}(x) \rightarrow 0$ a.e. because $\sum_n m(E_n) < \infty$ and so, by the Borel Cantelli lemma, there is a set of measure zero N such that for $x \notin N$, x is in only finitely many of the E_n . In particular, $\mathcal{R}_{E_n}(x) = 0$ for all n large enough if $x \notin N$. Then by the dominated convergence theorem, the inequality 11.22 cannot be valid for all n because the limit of the integral on the left equals 0. This is a contradiction. Hence F must be absolutely continuous after all.

Next suppose the function F is absolutely continuous. Suppose $m(E) = 0$. Does it follow that $\mu(E) = 0$? Let $\varepsilon > 0$ be given. Let δ correspond to $\varepsilon/2$ in the definition of absolute continuity. Let $E \subseteq V$ where V is an open set such that $m(V) < \delta$. By Theorem 3.11.8, $V = \cup_i (a_i, b_i)$ where these open intervals are disjoint. It follows that for each n , $\frac{\varepsilon}{2} > \sum_{i=1}^n F(b_i) - F(a_i) = \mu(\cup_{i=1}^n (a_i, b_i))$. Then letting $n \rightarrow \infty$, $\varepsilon > \frac{\varepsilon}{2} \geq \mu(\cup_{i=1}^{\infty} (a_i, b_i)) = \mu(V) \geq \mu(E)$. Since $\varepsilon > 0$ is arbitrary, it follows that $\mu(E) = 0$ and so $\mu \ll m$. ■

An example which shows that increasing and continuous is not enough, see Problem 5 on Page 269.

Corollary 11.14.4 F is increasing on $[a, b]$ and absolutely continuous if and only if $F'(x)$ exists for a.e. x and F' is in $L^1([a, b], m)$ and for every x, y such that $a \leq x \leq y \leq b$

$$F(y) - F(x) = \int_x^y F'(t) dm$$

Proof: \Rightarrow Suppose first that F is absolutely continuous. Then by Theorem 11.13.2, for μ defined above, $\mu(E) = \int_E D_m \mu(x) dm$ for all E Borel. In particular,

$$F(y) - F(x) = \int_{(x,y)} D_m \mu(t) dm(t) \quad (11.23)$$

Since $D_m \mu$ is in $L^1([a, b], m)$, it follows that almost every point is a Lebesgue point and so for such Lebesgue points x ,

$$\begin{aligned} \left| \frac{F(x+h) - F(x)}{h} - D_m \mu(x) \right| &= \left| \frac{1}{h} \int_{[x, x+h]} (D_m \mu(t) - D_m \mu(x)) dm(t) \right| \\ &\leq 2 \left| \frac{1}{2h} \int_{[x-h, x+h]} |D_m \mu(t) - D_m \mu(x)| dm(t) \right| \end{aligned}$$

which converges to 0 as $h \rightarrow 0$ since x is a Lebesgue point. Similarly, at each Lebesgue point, $\lim_{h \rightarrow 0} \frac{F(x) - F(x-h)}{h} = D_m \mu(x)$. Thus F is differentiable at each Lebesgue point and the derivative equals $D_m \mu$ at these points. Now 11.23 yields the desired result that the function can be recovered from integrating its derivative.

\Leftarrow Next suppose $F(y) - F(x) = \int_x^y F'(t) dm$ where $F'(t)$ exists a.e. and F' is in L^1 . Then if $\{I_i\}_i$ are nonoverlapping intervals, $\int_{\cup_i I_i} F'(t) dm = m(\cup_i F(I_i)) < \varepsilon$ if $m(\cup_i I_i)$ is small enough, as an application of the dominated convergence theorem or as in the first part of Theorem 11.14.3. ■

The importance of the intervals being non overlapping is discussed in the following proposition. I think it is also interesting that it implies F is Lipschitz. In this proposition, F is defined on some interval, possibly \mathbb{R} .

Proposition 11.14.5 *Suppose Definition 11.14.2 is unchanged except do not insist that the intervals be non-overlapping. Then F is not just absolutely continuous, it is also Lipschitz. The converse also holds.*

Proof: Choose m such that $\sum_{n=m}^{\infty} 2^{-n} < \delta$ where δ goes with 1 in the definition of absolute continuity. Let $r \leq 1/2$ and this implies that any such choice of r yields $\sum_{n=m}^{\infty} r^n < \delta$. Let $E_{nr} \equiv [F' > r^{-n}]$. If any E_{nr} has measure zero, then F' is bounded off a set of measure zero and this is what is desired. Otherwise, each E_{nr} has a point of density. Let $h_n \in [r^{n-1}, r^{n-2})$ so $r^{-n+1}h_n \geq 1$ and let r be small enough and m large enough that $\sum_{n=m}^{\infty} h_n < \delta$. For t_{N+m} a point of density for $E_{(N+m)r}$ let $I_n \equiv (t_{N+m} - rh_n/2, t_{N+m} + rh_n/2)$. Pick N large, say $N > 9$. Make r smaller if necessary so that for $n \in [m, N+m]$, $\frac{m(E_{nr} \cap I_n)}{m(I_n)} > \frac{1}{2}$. Note that $E_{(N+m)r} \subseteq E_{nr}$ for $n < N+m$ and $t_{N+m} \in E_{nr}$ must be a point of density for E_{nr} . Then, since $F' > r^{-n}$ on E_{nr} , one obtains the following sequence of inequalities.

$$\begin{aligned} 4.5 &< \frac{N+1}{2} = \sum_{n=m}^{m+N} \frac{1}{2} < \\ &\sum_{n=m}^{m+N} \frac{\overbrace{r^{-n}rh_n}^{\geq 1} m(E_n \cap I_n)}{m(I_n)} = \sum_{n=m}^{m+N} r^{-n} m(I_n) \frac{m(E_n \cap I_n)}{m(I_n)} = \sum_{n=m}^{m+N} r^{-n} m(E_n \cap I_n) \\ &\leq \sum_{n=m}^{m+N} \int_{t_{N+m}-rh_n/2}^{t_{N+m}+rh_n/2} \mathcal{X}_{E_n} F'(t) dm \leq \sum_{n=m}^{m+N} F(t_{N+m} + rh_n/2) - F(t_{N+m} - rh_n/2) < 1 \end{aligned}$$

by assumption, since the sum of the lengths of the intervals is smaller than δ . Thus $4.5 \leq 1$, a contradiction. Hence some E_{nr} has measure zero and so F' is bounded by a constant K off a set of measure zero. Hence, $|F(s) - F(t)| = \left| \int_s^t F'(u) du \right| \leq K|s - t|$ showing that F is Lipschitz.

The other direction is fairly obvious. If F is Lipschitz continuous, with Lipschitz constant K then if $\sum_{i=1}^m |x_{i+1} - x_i| < \delta$, then $\sum_{i=1}^m |F(x_{i+1}) - F(x_i)| \leq K \sum_{i=1}^m |x_{i+1} - x_i|$ so if $\varepsilon > 0$ is given, let $\delta = \varepsilon/K$. ■

Note that when μ is the Lebesgue Stieltjes measure coming from increasing continuous F , it follows from the definition that $D_m \mu(x) = \lim_{h \rightarrow 0} \frac{F(x+h) - F(x-h)}{2h}$. The following is another characterization of absolute continuity.

Corollary 11.14.6 *Let μ be the Lebesgue Stieltjes measure described above for increasing F defined on \mathcal{F}_μ containing the Borel sets. Let $I \equiv \{x : D_m \mu(x) = \infty\}$. Then $\mu \ll m$ on $\mathcal{F}_m \cap \mathcal{F}_\mu$ if and only if $\mu(I)$ has measure 0. Here \mathcal{F}_m is the σ algebra of Lebesgue measurable sets.*

Proof: \Rightarrow If $\mu \ll m$, then by Theorem 11.13.2, for all $E \in \mathcal{F}_\mu \cap \mathcal{F}_m$, it follows that $\mu(E) = \int_E D_m \mu(x) dm$. Then by the fundamental theorem of calculus, Theorem 11.13.2, there is a set of m measure zero N such that off this set, $D_m \mu(x)$ exists and is in \mathbb{R} . Thus $N \supseteq I$ and by absolute continuity, $\mu(I) = 0$.

\Leftarrow Next suppose $\mu(I) = 0$. Then F has no jumps because if it did, then μ (a jump) > 0 and the jump is also contained in I . Let $m(E) = 0$ for E a bounded set. Then define

$$G_n \equiv \left\{ t : \liminf_{r \rightarrow 0} \frac{\mu(B(t, r))}{2r} \leq n \right\}, n \in \mathbb{N}$$

and note that $I = \left\{ t : \liminf_{r \rightarrow 0} \frac{\mu(B(t,r))}{2r} = \infty \right\}$. Consider $G_n \cap E$ where $m(E) = 0$. Let V be a bounded open set and $V \supseteq G_n \cap E$ and $m(V) < \varepsilon/n$. Then there is a Vitali cover of $G_n \cap E$ consisting of closed balls B such that $\mu(B) \leq nm(B)$, each ball contained in V . Then by the covering theorem Theorem 9.12.2 on Page 263, there is a disjoint union of these which covers $G_n \cap E$ called B_i , $\mu(E \cap G_n \setminus \cup_i B_i) = 0$ and so $\mu(E \cap G_n) = \sum_i \mu(B_i) \leq n \sum_i m(B_i) < nm(V) < \varepsilon$ and since ε is arbitrary, $\mu(E \cap G_n) = 0$. Now, since $\mu(I) = 0$,

$$\mu(E) = \mu(E \cap I) \cup \cup_n \mu(E \cap G_n) = \cup_n \mu(E \cap G_n) = \lim_{n \rightarrow \infty} \mu(E \cap G_n) = 0$$

In general, if $m(E) = 0$, let $E_n \equiv E \cap B(0, n)$. Then from what was just shown, $\mu(E_n) = 0$ and so, taking a limit, $\mu(E) = 0$ also. Thus $\mu \ll m$. ■

11.15 Total Variation

The total variation function of an absolutely continuous function is itself absolutely continuous. This is shown here along with some of its implications.

Definition 11.15.1 A finite subset, P of $[a, b]$ is called a partition of $[x, y] \subseteq [a, b]$ if $P = \{x_0, x_1, \dots, x_n\}$ where $x = x_0 < x_1 < \dots < x_n = y$. For $f : [a, b] \rightarrow \mathbb{R}$ and $P = \{x_0, x_1, \dots, x_n\}$ define $V_P[x, y] \equiv \sum_{i=1}^n |f(x_i) - f(x_{i-1})|$. Denoting by $\mathcal{P}[x, y]$ the set of all partitions of $[x, y]$ define $V[x, y] \equiv \sup_{P \in \mathcal{P}[x, y]} V_P[x, y]$. For simplicity, $V[a, x]$ will be denoted by $V(x)$. It is called the total variation of the function f .

There are some simple facts about the total variation of an absolutely continuous function f which are contained in the next lemma.

Lemma 11.15.2 Let f be an absolutely continuous function defined on $[a, b]$ and let V be its total variation function as described above. Then V is an increasing bounded function. Also if P and Q are two partitions of $[x, y]$ with $P \subseteq Q$, then $V_P[x, y] \leq V_Q[x, y]$ and if $[x, y] \subseteq [z, w]$, $V[x, y] \leq V[z, w]$. If $P = \{x_0, x_1, \dots, x_n\}$ is a partition of $[x, y]$, then

$$V[x, y] = \sum_{i=1}^n V[x_i, x_{i-1}]. \quad (11.24)$$

Also if $y > x$, $V(y) - V(x) \geq |f(y) - f(x)|$ and the function, $x \rightarrow V(x) - f(x)$ is increasing. The total variation function V is absolutely continuous.

Proof: The claim that V is increasing is obvious as is the next claim about $P \subseteq Q$ leading to $V_P[x, y] \leq V_Q[x, y]$. To verify this, simply add in one point at a time and verify that from the triangle inequality, the sum involved gets no smaller. The claim that V is increasing consistent with set inclusion of intervals is also clearly true and follows directly from the definition.

Now let $t < V[x, y]$ where $P_0 = \{x_0, x_1, \dots, x_n\}$ is a partition of $[x, y]$. There exists a partition, P of $[x, y]$ such that $t < V_P[x, y]$. Without loss of generality it can be assumed that $\{x_0, x_1, \dots, x_n\} \subseteq P$ since if not, you can simply add in the points of P_0 and the resulting sum for the total variation will get no smaller. Let P_i be those points of P which are contained in $[x_{i-1}, x_i]$. Then $t < V_P[x, y] = \sum_{i=1}^n V_{P_i}[x_{i-1}, x_i] \leq \sum_{i=1}^n V[x_{i-1}, x_i]$. Since $t < V[x, y]$ is arbitrary,

$$V[x, y] \leq \sum_{i=1}^n V[x_i, x_{i-1}] \quad (11.25)$$

Note that 11.25 does not depend on f being absolutely continuous.

Suppose now that f is absolutely continuous. Let δ correspond to $\varepsilon = 1$. Then if $[x, y]$ is an interval of length no larger than δ , the definition of absolute continuity implies $V[x, y] < 1$. Then from 11.25, $V[a, n\delta] \leq \sum_{i=1}^n V[a + (i-1)\delta, a + i\delta] < \sum_{i=1}^n 1 = n$. Thus V is bounded on $[a, b]$. Now let P_i be a partition of $[x_{i-1}, x_i]$ such that $V_{P_i}[x_{i-1}, x_i] > V[x_{i-1}, x_i] - \frac{\varepsilon}{n}$. Then letting $P = \cup P_i$,

$$-\varepsilon + \sum_{i=1}^n V[x_{i-1}, x_i] < \sum_{i=1}^n V_{P_i}[x_{i-1}, x_i] = V_P[x, y] \leq V[x, y].$$

Since ε is arbitrary, 11.24 follows from this and 11.25.

Now let $x < y$. $V(y) - f(y) - (V(x) - f(x)) =$

$$V(y) - V(x) - (f(y) - f(x)) \geq V(y) - V(x) - |f(y) - f(x)| \geq 0.$$

It only remains to verify that V is absolutely continuous.

Let $\varepsilon > 0$ be given and let δ correspond to $\varepsilon/2$ in the definition of absolute continuity applied to f . Suppose $\sum_{i=1}^n |y_i - x_i| < \delta$ and consider $\sum_{i=1}^n |V(y_i) - V(x_i)|$. By 11.25 this last is no larger than $\sum_{i=1}^n V[x_i, y_i]$. Now let P_i be a partition of $[x_i, y_i]$ such that $V_{P_i}[x_i, y_i] + \frac{\varepsilon}{2n} > V[x_i, y_i]$. Then by the definition of absolute continuity,

$$\sum_{i=1}^n |V(y_i) - V(x_i)| = \sum_{i=1}^n V[x_i, y_i] \leq \sum_{i=1}^n V_{P_i}[x_i, y_i] + \eta < \varepsilon/2 + \varepsilon/2 = \varepsilon$$

and shows V is absolutely continuous as claimed. ■

Now with the above results, the following is the main result on absolutely continuous functions.

Theorem 11.15.3 *Let $f: [a, b] \rightarrow \mathbb{R}$ be a function. Then f is absolutely continuous if and only if $f'(t)$ exists a.e., f' is in $L^1([a, b], m)$, and for every $a \leq x \leq y \leq b$,*

$$f(y) - f(x) = \int_x^y f'(t) dt \equiv \int_{[x, y]} f'(t) dm(t)$$

Proof: Suppose f is absolutely continuous. Using Lemma 11.15.2, $f(x) = V(x) - (V(x) - f(x))$, the difference of two increasing functions, both of which are absolutely continuous. See Problem 1 on Page 349. Denote the derivatives of these two increasing functions by k and l respectively. Then for $x \leq y$,

$$f(y) - f(x) = \int_{[x, y]} k(t) dm(t) - \int_{[x, y]} l(t) dm(t)$$

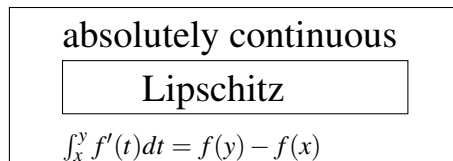
Letting $g(t) \equiv k(t) - l(t)$, it follows that $f(y) - f(x) = \int_x^y g(t) dt$ where $g \in L^1$. Then from the fundamental theorem of calculus, Theorem 11.4.2, if x is a Lebesgue point of g , not equal to one of the end points.

$$\left| \frac{f(x+h) - f(x)}{h} - g(x) \right| = \left| \frac{1}{h} \int_x^{x+h} g(t) - g(x) dt \right| \leq 2 \left(\frac{1}{2h} \int_{x-h}^{x+h} |g(t) - g(x)| dt \right)$$

which converges to 0. Hence $g(x) = f'(x)$ a.e.

Next suppose $f(y) - f(x) = \int_x^y f'(t) dt$ where $f' \in L^1$. If f is not absolutely continuous, there exists $\varepsilon > 0$ and sets V_n each of which is the union of non-overlapping intervals such that $m(V_n) < 2^{-n}$ but $\int_{V_n} |f'(t)| dt \geq \varepsilon$. However, by the Borel Cantelli lemma, there exists a set of measure zero N such that for $x \notin N$, it follows that x is in only finitely many of the V_n . Thus $\mathcal{R}_{V_n}(x) \rightarrow 0$. Then a use of the dominated convergence theorem implies that $\lim_{n \rightarrow \infty} \int_{V_n} |f'(t)| dt = 0$ which is a contradiction. Thus f must be absolutely continuous. ■

The following picture illustrates the main items shown so far about functions of one variable.



11.16 Exercises

1. Show that if f is absolutely continuous on $[a, b]$ and if $V(x)$ is the total variation of f on $[0, x]$, then V is also absolutely continuous.
2. In Problem 5 on Page 310, you showed that if $f \in L^1(\mathbb{R}^p)$, there exists h which is continuous and equal to 0 off some compact set such that $\int |f - h| dm < \varepsilon$. Define $f_y(x) \equiv f(x - y)$. Explain why f_y is Lebesgue measurable and $\int |f_y| dm_p = \int |f| dm_p$. Now justify the following formula. $\int |f_y - f| dm_p \leq \int |f_y - h_y| dm_p + \int |h_y - h| dm_p + \int |h - f| dm_p \leq 2\varepsilon + \int |h_y - h| dm_p$. Now explain why the last term is less than ε if $\|y\|$ is small enough. Explain continuity of translation in $L^1(\mathbb{R}^p)$ which says that $\lim_{y \rightarrow 0} \int |f_y - f| dm_p = 0$.
3. This problem will help to understand that a certain kind of function exists. Let $f(x) = e^{-1/x^2}$ if $x \neq 0$ and let $f(x) = 0$ if $x = 0$. Show that f is infinitely differentiable. Note that you only need to be concerned with what happens at 0. There is no question elsewhere. This is a little fussy but is not too hard.
4. †Let $f(x)$ be as given above. Now let $\hat{f}(x) = f(x)$ if $x \leq 0$ and let $\hat{f}(x) = 0$ if $x > 0$. Show that $\hat{f}(x)$ is also infinitely differentiable. Now let $r > 0$ and define $g(x) \equiv \hat{f}(-(x-r))\hat{f}(x+r)$. Show that g is infinitely differentiable and vanishes for $|x| \geq r$. Let $\psi(x) = \prod_{k=1}^p g(x_k)$. For $U = B(0, 2r)$ with the norm given by $\|x\| = \max\{|x_k|, k \leq p\}$, show that $\psi \in C_c^\infty(U)$.
5. †Using the above problem, show there exists $\psi \geq 0$ such that $\psi \in C_c^\infty(B(0, 1))$ and $\int \psi dm_p = 1$. Now define $\psi_n(x) \equiv n^p \psi(nx)$. Show that ψ_n equals zero off a compact subset of $B(0, \frac{1}{n})$ and $\int \psi_n dm_p = 1$. We say that $\text{spt}(\psi_n) \subseteq B(0, \frac{1}{n})$. $\text{spt}(f)$ is defined as the closure of the set on which f is not equal to 0. Such a sequence of functions as just defined $\{\psi_n\}$ where $\int \psi_n dm_p = 1$ and $\psi_n \geq 0$ and $\text{spt}(\psi_n) \subseteq B(0, \frac{1}{n})$ is called a **mollifier**.
6. †It is important to be able to approximate functions with those which are infinitely differentiable. Suppose $f \in L^1(\mathbb{R}^p)$ and let $\{\psi_n\}$ be a mollifier as above. We define the convolution as follows. $f * \psi_n(x) \equiv \int f(x - y) \psi_n(y) dm_p(y)$ Here the notation means that the variable of integration is y . Show that $f * \psi_n(x)$ exists

and equals $\int \psi_n(\mathbf{x} - \mathbf{y}) f(\mathbf{y}) dm_p(\mathbf{y})$. Now show using the dominated convergence theorem that $f * \psi_n$ is infinitely differentiable. Next show that

$$\lim_{n \rightarrow \infty} \int |f(\mathbf{x}) - f * \psi_n(\mathbf{x})| dm_p = 0.$$

Thus, in terms of being close in $L^1(\mathbb{R}^p)$, every function in $L^1(\mathbb{R}^p)$ is close to one which is infinitely differentiable.

7. \uparrow From Problem 5 above and $f \in L^1(\mathbb{R}^p)$, there exists $h \in C_c(\mathbb{R}^p)$, continuous and $\text{spt}(h)$ a compact set, such that $\int |f - h| dm_p < \varepsilon$. Now consider $h * \psi_n$. Show that this function is in $C_c^\infty(\text{spt}(h) + B(\mathbf{0}, \frac{2}{n}))$. The notation means you start with the compact set $\text{spt}(h)$ and fatten it up by adding the set $B(\mathbf{0}, \frac{1}{n})$. It means $\mathbf{x} + \mathbf{y}$ such that $\mathbf{x} \in \text{spt}(h)$ and $\mathbf{y} \in B(\mathbf{0}, \frac{1}{n})$. Show the following. For all n large enough, $\int |f - h * \psi_n| dm_p < \varepsilon$ so one can approximate with a function which is infinitely differentiable and also has compact support. Also show that $h * \psi_n$ converges uniformly to h . If h is a function in $C^k(\mathbb{R}^n)$ in addition to being continuous with compact support, show that for each $|\alpha| \leq k$, $D^\alpha(h * \psi_n) \rightarrow D^\alpha h$ uniformly. **Hint:** If you do this for a single partial derivative, you will see how it works in general.
8. \uparrow Let $f \in L^1(\mathbb{R})$. Show that $\lim_{n \rightarrow \infty} \int f(x) \sin(nx) dm = 0$. **Hint:** Use the result of the above problem to obtain $g \in C_c^\infty(\mathbb{R})$, continuous and zero off a compact set, such that $\int |f - g| dm < \varepsilon$. Then show that $\lim_{n \rightarrow \infty} \int g(x) \sin(nx) dm(x) = 0$. You can do this by integration by parts. Then consider this. $|\int f(x) \sin(nx) dm| =$

$$\begin{aligned} & \left| \int f(x) \sin(nx) dm - \int g(x) \sin(nx) dm \right| + \left| \int g(x) \sin(nx) dm \right| \\ & \leq \int |f - g| dm + \left| \int g(x) \sin(nx) dm \right| \end{aligned}$$

This is the celebrated Riemann Lebesgue lemma which is the basis for all theorems about pointwise convergence of Fourier series and Fourier integrals.

9. As another application of theory of regularity, here is a very important result. Suppose $f \in L^1(\mathbb{R}^p)$ and for every $\psi \in C_c^\infty(\mathbb{R}^p)$ $\int f \psi dm_p = 0$. Show that then it follows $f(x) = 0$ for a.e. x . That is, there is a set of measure zero such that off this set f equals 0. **Hint:** What you can do is to let E be a measurable set which is bounded and let $K_n \subseteq E \subseteq V_n$ where $m_p(V_n \setminus K_n) < 2^{-n}$. Here K_n is compact and V_n is open. By an earlier exercise, Problem 11 on Page 259, there exists a function ϕ_n which is continuous, has values in $[0, 1]$ equals 1 on K_n and $\text{spt}(\phi_n) \subseteq V$. To get this last part, show there exists W_n open such that $\bar{W}_n \subseteq V_n$ and W_n contains K_n . Then you use the problem to get $\text{spt}(\phi_n) \subseteq \bar{W}_n$. Now you form $\eta_n = \phi_n * \psi_l$ where $\{\psi_l\}$ is a mollifier. Show that for l large enough, η_n has values in $[0, 1]$, $\text{spt}(\eta_n) \subseteq V_n$ and $\eta_n \in C_c^\infty(V_n)$. Now explain why $\eta_n \rightarrow \mathcal{X}_E$ off a set of measure zero. To do this, you might want to consider the Borel Cantelli lemma, Lemma 9.2.5 on Page 243. Then

$$\left| \int f \mathcal{X}_E dm_p \right| = \left| \int f(\mathcal{X}_E - \eta_n) dm_p \right| + \left| \int f \eta_n dm_p \right| = \left| \int f(\mathcal{X}_E - \eta_n) dm_p \right|$$

Now explain why this converges to 0 on the right. This will involve the dominated convergence theorem. Conclude that $\int f \mathcal{X}_E dm_p = 0$ for every bounded measurable

set E . Show that this implies that $\int f \mathcal{X}_E dm_p = 0$ for every measurable E . Explain why this requires $f = 0$ a.e.

10. Suppose f, g are absolutely continuous on $[a, b]$. Prove the usual integration by parts formula. **Hint:** You might try the following:

$$\int_a^b f'(t) g(t) dm(t) = \int_a^b f'(t) \left(\int_a^t g'(s) dm(s) + g(a) \right) dm(t)$$

Now do one integration and then interchange the order of integration.

11. Let $x \rightarrow F(f(x))$ be absolutely continuous whenever $x \rightarrow f(x)$ is absolutely continuous. Then F is Lipschitz. This is due to G. M. Fishtenholz. **Hint:** Reduce to Proposition 11.14.5 using an appropriate Lipschitz continuous function f .
12. Suppose $g: [c, d] \rightarrow [a, b]$ and is absolutely continuous and increasing and $f: [a, b] \rightarrow \mathbb{R}$ is Lipschitz continuous. Show that then $f \circ g$ is absolutely continuous and

$$\int_a^b f(t) dm(t) = \int_c^d f'(g(s)) g'(s) ds = f(b) - f(a)$$

13. If $f \in L^1(\Omega, \mu)$, show that $\lim_{\mu(E) \rightarrow 0} \int_E |f| d\mu = 0$. Defining $F(x) \equiv \int_a^x f(t) dm(t)$ for $f \in L^1([a, b], m)$, verify that F is absolutely continuous with $F'(x) = f(x)$ a.e.
14. Show that if f is absolutely continuous, as defined in Definition 11.14.2, then it is of bounded variation.
15. Let $f: [a, b] \rightarrow \mathbb{R}$ be absolutely continuous. Show that in fact, the total variation of f on $[a, b]$ is $\int_a^b |f'| dm$. **Hint:** One direction is easy, that $V[a, b] \leq \int_a^b |f'| dm$. To do the other direction, show there is a sequence of **step functions** $s_n(t) \equiv \sum_{k=1}^{m_n} \alpha_k^n \mathcal{X}_{I_k^n}(t)$ which converges to $\text{sgn } f'$ pointwise a.e., $I_k^n = (c_{k-1}^n, c_k^n)$. This will involve regularity notions. Explain why it can be assumed each $|\alpha_k^n| \leq 1$. Then

$$\left| \int_a^b f' s_n \right| = \left| \sum_{k=1}^{m_n} \alpha_k^n \int_{I_k^n} f' \right| \leq \sum_{k=1}^{m_n} |f(c_k^n) - f(c_{k-1}^n)| \leq V([a, b], f)$$

Now pass to a limit using the dominated convergence theorem.

16. Let $F(x) = \left(\int_0^x e^{-t^2} dt \right)^2$. Justify the following:

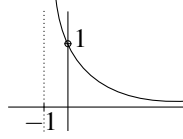
$$F'(x) = 2 \left(\int_0^x e^{-t^2} dt \right) e^{-x^2} = 2xe^{-x^2} \left(\int_0^1 e^{-x^2 t^2} dt \right) = 2x \left(\int_0^1 e^{-x^2(t^2+1)} dt \right)$$

Now integrate.

$$\begin{aligned} F(x) &= \int_0^x \int_0^1 2ue^{-u^2(t^2+1)} dt du = \int_0^1 \int_0^x 2ue^{-u^2(t^2+1)} du dt \\ &= \int_0^1 -e^{-u^2(t^2+1)} \frac{1}{1+t^2} \Big|_0^x dt = \int_0^1 \left(\frac{1}{1+t^2} - e^{-x^2} \frac{1}{1+t^2} \right) dt \end{aligned}$$

Now let $x \rightarrow \infty$ and conclude $F(\infty) = \left(\int_0^\infty e^{-t^2} dt \right)^2 = \int_0^1 \frac{1}{1+t^2} dt = \frac{\pi}{4}$.

17. This problem outlines an approach to Stirling's formula following [49] and [7]. From the above problems, $\Gamma(n+1) = n!$ for $n \geq 0$. Consider more generally $\Gamma(x+1)$ for $x > 0$. Actually, we will always assume $x > 1$ since it is the limit as $x \rightarrow \infty$ which is of interest. $\Gamma(x+1) = \int_0^\infty e^{-t} t^x dt$. Change variables letting $t = x(1+u)$ to obtain $\Gamma(x+1) = x^{x+1} e^{-x} \int_{-1}^\infty ((1+u)e^{-u})^x du$. Next let $h(u)$ be such that $h(0) = 1$ and $(1+u)e^{-u} = \exp\left(-\frac{u^2}{2}h(u)\right)$. Show that the thing which works is $h(u) = \frac{2}{u^2}(u - \ln(1+u))$. Use L'Hospital's rule to verify that the limit of $h(u)$ as $u \rightarrow 0$ is 1. The graph of h is illustrated in the following picture. Verify that its graph is like this, with an asymptote at $u = -1$ decreasing and equal to 1 at 0 and converging to 0 as $u \rightarrow \infty$.



Next change the variables again letting $u = s\sqrt{\frac{2}{x}}$. This yields, from the original description of h

$$\Gamma(x+1) = x^x e^{-x} \sqrt{2x} \int_{-\sqrt{x/2}}^\infty \exp\left(-s^2 h\left(s\sqrt{\frac{2}{x}}\right)\right) ds$$

For $s < 1$, $h\left(s\sqrt{\frac{2}{x}}\right) > 2 - 2\ln 2 = 0.61371$ so the above integrand is dominated by $e^{-(2-2\ln 2)s^2}$. Consider the integrand in the above for $s > 1$. Show that the exponent part is

$$-\left(\sqrt{2}\sqrt{x}s - x \ln\left(1 + s\sqrt{\frac{2}{x}}\right)\right)$$

The expression $\left(\sqrt{2}\sqrt{x}s - x \ln\left(1 + s\sqrt{\frac{2}{x}}\right)\right)$ is increasing in x . You can show this by fixing s and taking a derivative with respect to x . Therefore, it is larger than $\sqrt{2}\sqrt{1}s - \ln\left(1 + s\sqrt{\frac{2}{1}}\right)$ and so

$$\begin{aligned} \exp\left(-s^2 h\left(s\sqrt{\frac{2}{x}}\right)\right) &\leq \exp\left(-\left(\sqrt{2}\sqrt{1}s - \ln\left(1 + s\sqrt{\frac{2}{1}}\right)\right)\right) \\ &= (1 + s\sqrt{2}) e^{-\sqrt{2}s} \end{aligned}$$

Thus, there exists a dominating function for $\mathcal{X}_{[-\sqrt{x/2}, \infty]}(s) \exp\left(-s^2 h\left(s\sqrt{\frac{2}{x}}\right)\right)$ and these functions converge pointwise to $\exp(-s^2)$ as $x \rightarrow \infty$ so by the dominated convergence theorem,

$$\lim_{x \rightarrow \infty} \int_{-\sqrt{x/2}}^\infty \exp\left(-s^2 h\left(s\sqrt{\frac{2}{x}}\right)\right) ds = \int_{-\infty}^\infty e^{-s^2} ds = \sqrt{\pi}$$

See Problem 16. This yields a general Stirling's formula, $\lim_{x \rightarrow \infty} \frac{\Gamma(x+1)}{x^x e^{-x} \sqrt{2x}} = \sqrt{\pi}$.

18. Let $\bar{\mu}$ be the completion of the measure $\mu = \prod_{i=1}^p \mu_i$ on the product σ algebra $\prod_{i=1}^p \mathcal{F}_i$. Show that if $f \geq 0$ is $\overline{\prod_{i=1}^p \mathcal{F}_i}$ measurable, then there are two functions g, h which are $\prod_{i=1}^p \mathcal{F}_i$ measurable and $g \geq f \geq h$ while $\mu([g - h > 0]) = 0$. Assume the original measure spaces are finite or σ finite.
19. In the representation theorem for positive linear functionals, show that if $\mu(V) = \sup\{Lf : f \prec V\}$, then the σ algebra and measure representing the functional are unique.
20. Show that $\int_0^\infty e^{-x^2} dx = \frac{1}{2}\sqrt{\pi}$. **Hint:** First verify the integral is finite. You might use monotone convergence theorem to do this. It is easier than the stuff you worried about in beginning calculus. Next let $I = \int_0^\infty e^{-x^2} dx$ so $I^2 = \int_0^\infty \int_0^\infty e^{-(x^2+y^2)} dx dy$. Now change the variables using polar coordinates. It is all justified by the big change of variables theorem we have done. This becomes an easy problem when you do this.
21. Show that $\int_{-\infty}^\infty \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = 1$. Here σ is a positive number called the standard deviation and μ is a number called the mean. **Hint:** Just use the above result to find $\int_{-\infty}^\infty e^{-x^2} dx$ and then change the variables in this one.
22. The Gamma function is $\Gamma(\alpha) \equiv \int_0^\infty e^{-t} t^{\alpha-1} dt$. Verify that this function is well defined in the sense that it is finite for all $\alpha > 0$. Next verify that $\Gamma(1) = 1 = \Gamma(2)$ and $\Gamma(x+1) = \Gamma(x)x$. Conclude that $\Gamma(n+1) = n!$ if n is an integer. Now consider $\Gamma(1/2)$. This is $\int_0^\infty e^{-t} t^{-1/2} dt$. Change the variables in this integral. You might let $t = u^2$. Then consider the above problem.
23. Let $p, q > 0$ and define $B(p, q) = \int_0^1 x^{p-1} (1-x)^{q-1} dx$. Show that the following identity holds: $\Gamma(p)\Gamma(q) = B(p, q)\Gamma(p+q)$. **Hint:** It is fairly routine if you start with the left side and proceed to change variables.
24. Let E be a Lebesgue measurable set in \mathbb{R} . Suppose $m(E) > 0$. Consider the set $E - E = \{x - y : x \in E, y \in E\}$. Show that $E - E$ contains an interval. This is an amazing result. Recall the case of the fat Cantor set which contained no intervals but had positive measure. **Hint:** Without loss of generality, you can assume E is bounded. Let $f(x) = \int \mathcal{X}_E(t) \mathcal{X}_E(x+t) dt$. Explain why f is continuous at 0 and $f(0) > 0$ and use continuity of translation in L^1 . To see it is continuous,

$$\begin{aligned} |f(x) - f(\hat{x})| &\leq \int \mathcal{X}_E(t) |\mathcal{X}_E(x+t) - \mathcal{X}_E(\hat{x}+t)| dt \\ &\leq \int |\mathcal{X}_E(x+t) - \mathcal{X}_E(\hat{x}+t)| dt \end{aligned}$$

Now explain why this is small whenever $\hat{x} - x$ is small due to continuity of translation in $L^1(\mathbb{R})$. Thus $f(0) = m(E) > 0$ and so by continuity, $f > 0$ near 0. If the integral is nonzero, what can you say about the integrand? You must have for all $x \in (-\delta, \delta)$ both $x+t \in E$ and $t \in E$. Now consider this a little.

25. Does there exist a closed uncountable set which is contained in the set of irrational numbers? If so, explain why and if not, explain why. Thus this uncountable set has no rational number as a limit point.

26. Find the area of the bounded region R , determined by $5x + y = 1$, $5x + y = 9$, $y = 2x$, and $y = 5x$.
27. Here are three vectors. $(1, 2, 3)^T$, $(1, 0, 1)^T$, and $(2, 1, 0)^T$. These vectors determine a parallelepiped, R , which is occupied by a solid having density $\rho = y$. Find the mass of this solid. To find the mass of the solid, you integrate the density. Thus, if P is this parallelepiped, the mass is $\int_P y dm_3$. **Hint:** Let $h : [0, 1]^3 \rightarrow P$ be given by $h(t_1, t_2, t_3) = t_1 \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} + t_2 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + t_3 \begin{pmatrix} 2 \\ 1 \\ 0 \end{pmatrix}$ then by definition of what is meant by a parallelepiped, $h([0, 1]^3) = P$ and h is one to one and onto.
28. Suppose $f, g \in L^1(\mathbb{R}^p)$. Define $f * g(x)$ by $\int f(x - y)g(y) dm_p(y)$. First show using the preceding problem that there is no loss of generality in assuming that both f, g are Borel measurable. Next show this makes sense a.e. x and that in fact for a.e. x , $\int |f(x - y)||g(y)| dm_p(y) < \infty$. Next show $\int |f * g(x)| dm_p(x) \leq \int |f| dm_p \int |g| dm_p$. **Hint:** You can use Fubini's theorem to write

$$\begin{aligned} \int \int |f(x - y)||g(y)| dm_p(y) dm_p(x) &= \\ \int \int |f(x - y)||g(y)| dm_p(x) dm_p(y) &= \int |f(z)| dm_p \int |g(y)| dm_p. \end{aligned}$$

29. Suppose $X : (\Omega, \mathcal{F}, P)$ where P is a probability measure and suppose $X : \Omega \rightarrow \mathbb{R}$ is measurable. That is, $X^{-1}(\text{open set}) \in \mathcal{F}$. Then consider the distribution measure λ_X defined on the Borel sets of \mathbb{R}^p and given as follows. $\lambda_X(E) = P(X \in E)$. Explain why this is a probability measure on $\mathcal{B}(\mathbb{R})$ and why $X^{-1}(B) \in \mathcal{F}$ whenever B is a Borel set. Next show that if $X \in L^1(\Omega)$, $\int_{\Omega} X dP = \int_{\mathbb{R}} x d\lambda_X$. Suppose h is a complex valued Borel measurable function defined on \mathbb{R} which is also bounded. Show that

$$\int h(x) d\lambda_X(x) = \int h(X(\omega)) dP$$

Hint: Recall that from the definition of the integral,

$$\int_{\mathbb{R}} |x| d\lambda_X = \int_0^{\infty} \lambda_X(|x| > \alpha) d\alpha = \int_0^{\infty} P(|X| > \alpha) d\alpha = \int_{\Omega} |X| dP < \infty$$

30. Let $h : U \rightarrow h(U)$ be one to one and C^1 . Use the inverse function theorem to give a much easier proof of the change of variables formula.
31. If a continuous function is one to one on a compact set, explain why its inverse is continuous.
32. Suppose U is a nonempty set in \mathbb{R}^p . Let ∂U consist of the points $p \in \mathbb{R}^p$ such that $B(p, r)$ contains points of U as well as points of $\mathbb{R}^p \setminus U$. Show that U is contained in the union of the interior of U , denoted as $\text{int}(U)$ with ∂U . Now suppose that $f : U \rightarrow \mathbb{R}^p$ and is one to one and continuous. Explain why $\text{int}(f(U))$ equals $f(\text{int}(U))$.
33. Prove the Radon Nikodym theorem, Theorem 11.13.2 in case $\lambda \ll \mu$ another way by using the earlier general Radon Nikodym theorem, Theorem 10.13.7 or its corollary and then identifying the function ensured by that theorem with the symmetric derivative, using the fundamental theorem of calculus, Theorem 11.4.2.

34. Suppose \mathbf{x} is a Lebesgue point of f with respect to Lebesgue measure so that

$$\lim_{r \rightarrow 0} \frac{1}{m_p(\mathbf{x}, r)} \int_{B(\mathbf{x}, r)} |f(\mathbf{x}) - f(\mathbf{y})| dm_p(y) = 0.$$

Suppose $\mathbf{x} \in E_r \subseteq B(\mathbf{x}, \sigma r)$, E_r a measurable set and there exists δ such that

$$\delta m_p(B(\mathbf{x}, r)) < m_p(E_r)$$

for all r . Verify that $\lim_{r \rightarrow 0} \frac{1}{m_p(E_r)} \int_{E_r} |f(\mathbf{x}) - f(\mathbf{y})| dm_p(y) = 0$.

35. Generalize Corollary 11.14.6 to the case where m is replaced by m_p and μ is some Radon measure on \mathcal{F}_μ a σ algebra of sets of \mathbb{R}^p or an open subset of \mathbb{R}^p .
36. If F is increasing and continuous and I consists of all t such that $F'_+(t) = \infty$ or $F'_-(t) = \infty$, show that if $\mu(I) = 0$ for μ the Lebesgue Stieltjes measure associated with F , then F is absolutely continuous.
37. Let $I \equiv \left\{ t : \min \left(\liminf_{h \rightarrow 0} \frac{F(t+h) - F(t)}{h}, \liminf_{h \rightarrow 0} \frac{F(t) - F(t-h)}{h} \right) = \infty \right\}$ where F is increasing and continuous. Let μ be the Lebesgue Stieltjes measure coming from F . Show that if $\mu(I) = 0$, then $\mu \ll m$. This I is the set where $F'(t) = \infty$. Conversely, if $\mu \ll m$, then F is continuous and $\mu(I) = \infty$. **Hint:** Let

$$G_n \equiv \left\{ t : \min \left(\liminf_{h \rightarrow 0} \frac{F(t+h) - F(t)}{h}, \liminf_{h \rightarrow 0} \frac{F(t) - F(t-h)}{h} \right) \leq n \right\}$$

and follow the idea of Corollary 11.14.6 using a covering theorem, Theorem 9.12.2. The assumption that F is continuous is needed to say that, for example, $\frac{F(t+h) - F(t)}{h} = \frac{\mu([t, t+h])}{h}$.

38. Let $\mathbf{h} : U \rightarrow \mathbf{h}(U)$ be C^1 and one to one. Give a much easier change of variables formula using the covering theorem for Vitali covers and the material on linear mappings. Then extend to the case where \mathbf{h} is maybe not one to one using the inverse function theorem. You might first prove such a theorem for f continuous with compact support and use the Riesz representation theorem for positive linear functionals.

Chapter 12

The L^p Spaces

12.1 Basic Inequalities and Properties

One of the main applications of the Lebesgue integral is to the study of various sorts of functions space. These are vector spaces whose elements are functions of various types. One of the most important examples of a function space is the space of measurable functions whose absolute values are p^{th} power integrable where $p \geq 1$. These spaces, referred to as L^p spaces, are very useful in applications. In the chapter $(\Omega, \mathcal{S}, \mu)$ will be a measure space.

Definition 12.1.1 Let $1 \leq p < \infty$. Define

$$L^p(\Omega) \equiv \{f : f \text{ is measurable and } \int_{\Omega} |f(\omega)|^p d\mu < \infty\}$$

In terms of the distribution function,

$$L^p(\Omega) = \{f : f \text{ is measurable and } \int_0^{\infty} pt^{p-1} \mu(|f| > t) dt < \infty\}$$

because $\int_0^{\infty} \mu(|f|^p > t) dt = \int_0^{\infty} \mu(|f| > t^{1/p}) dt = \int_0^{\infty} ps^{p-1} \mu(|f| > s) ds$.

For each $p > 1$ define q by $\frac{1}{p} + \frac{1}{q} = 1$. Often one uses p' instead of q in this context.

$L^p(\Omega)$ is a vector space and has a norm. This is similar to the situation for \mathbb{R}^n but the proof requires the following fundamental inequality. When $p = 1$, we use the symbol ∞ to represent q . The space $L^{\infty}(\Omega)$ will be discussed later.

Theorem 12.1.2 (Holder's inequality) If f and g are measurable functions, then if $p > 1$,

$$\int |f| |g| d\mu \leq \left(\int |f|^p d\mu \right)^{\frac{1}{p}} \left(\int |g|^q d\mu \right)^{\frac{1}{q}}. \quad (12.1)$$

Proof: First recall Lemma 4.3.10, stated here for convenience.

Lemma 12.1.3 If $p > 1$, and $0 \leq a, b$ then $ab \leq \frac{a^p}{p} + \frac{b^q}{q}$. Equality occurs when $a^p = b^q$.

Proof of Holder's inequality: If either $\int |f|^p d\mu$ or $\int |g|^q d\mu$ equals ∞ , the inequality 12.1 is obviously valid because $\infty \geq$ anything. If either $\int |f|^p d\mu$ or $\int |g|^q d\mu$ equals 0, then $f = 0$ a.e. or $g = 0$ a.e. and so in this case the left side of the inequality equals 0 and so the inequality is therefore true. Therefore assume both $\int |f|^p d\mu$ and $\int |g|^q d\mu$ are less than ∞ and not equal to 0. Let $(\int |f|^p d\mu)^{1/p} = I(f)$ and let $(\int |g|^q d\mu)^{1/q} = I(g)$. Then using the lemma,

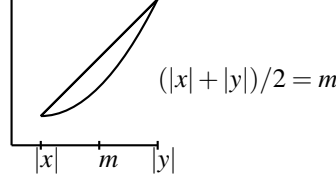
$$\int \frac{|f|}{I(f)} \frac{|g|}{I(g)} d\mu \leq \frac{1}{p} \int \frac{|f|^p}{I(f)^p} d\mu + \frac{1}{q} \int \frac{|g|^q}{I(g)^q} d\mu = 1.$$

Hence, $\int |f| |g| d\mu \leq I(f) I(g) = (\int |f|^p d\mu)^{1/p} (\int |g|^q d\mu)^{1/q}$. This proves Holder's inequality. ■

The following lemma will be needed.

Lemma 12.1.4 Suppose $x, y \in \mathbb{C}$. Then $|x + y|^p \leq 2^{p-1}(|x|^p + |y|^p)$.

Proof: The function $f(t) = t^p$ is concave up for $t \geq 0$ because $p > 1$. Therefore, the secant line joining two points on the graph of this function must lie above the graph of the function. This is illustrated in the following picture.



Since $\left(\frac{|x|+|y|}{2}\right)^p \leq \frac{|x|^p+|y|^p}{2}$, $|x + y|^p \leq (|x| + |y|)^p \leq 2^{p-1}(|x|^p + |y|^p)$ ■

Note that if $y = \phi(x)$ is any function for which the graph of ϕ is concave up, you could get a similar inequality by the same argument.

Corollary 12.1.5 (Minkowski inequality) Let $1 \leq p < \infty$. Then

$$\left(\int |f + g|^p d\mu\right)^{1/p} \leq \left(\int |f|^p d\mu\right)^{1/p} + \left(\int |g|^p d\mu\right)^{1/p}. \quad (12.2)$$

Proof: If $p = 1$, this is obvious. Let $p > 1$. Without loss of generality, assume $(\int |f|^p d\mu)^{1/p} + (\int |g|^p d\mu)^{1/p} < \infty$ and $(\int |f + g|^p d\mu)^{1/p} \neq 0$ or there is nothing to prove. Therefore, using the above lemma,

$$\int |f + g|^p d\mu \leq 2^{p-1} \left(\int |f|^p + |g|^p d\mu\right) < \infty.$$

Now $|f(\omega) + g(\omega)|^p \leq |f(\omega) + g(\omega)|^{p-1} (|f(\omega)| + |g(\omega)|)$. Also, it follows from the definition of p and q that $p - 1 = \frac{p}{q}$. Therefore, using this and Holder's inequality, $\int |f + g|^p d\mu \leq$

$$\begin{aligned} & \int |f + g|^{p-1} |f| d\mu + \int |f + g|^{p-1} |g| d\mu = \int |f + g|^{\frac{p}{q}} |f| d\mu + \int |f + g|^{\frac{p}{q}} |g| d\mu \\ & \leq \left(\int |f + g|^p d\mu\right)^{\frac{1}{q}} \left(\int |f|^p d\mu\right)^{\frac{1}{p}} + \left(\int |f + g|^p d\mu\right)^{\frac{1}{q}} \left(\int |g|^p d\mu\right)^{\frac{1}{p}}. \end{aligned}$$

Dividing both sides by $(\int |f + g|^p d\mu)^{\frac{1}{q}}$ yields 12.2. ■

The above theorem implies the following corollary.

Corollary 12.1.6 Let $f_i \in L^p(\Omega)$ for $i = 1, 2, \dots, n$. Then

$$\left(\int \left|\sum_{i=1}^n f_i\right|^p d\mu\right)^{1/p} \leq \sum_{i=1}^n \left(\int |f_i|^p d\mu\right)^{1/p}.$$

This shows that if $f, g \in L^p$, then $f + g \in L^p$. Also, it is clear that if a is a constant and $f \in L^p$, then $af \in L^p$ because $\int |af|^p d\mu = |a|^p \int |f|^p d\mu < \infty$. Thus L^p is a vector space and

a.) $(\int |f|^p d\mu)^{1/p} \geq 0$, $(\int |f|^p d\mu)^{1/p} = 0$ if and only if $f = 0$ a.e.

b.) $(\int |af|^p d\mu)^{1/p} = |a| (\int |f|^p d\mu)^{1/p}$ if a is a scalar.

c.) $(\int |f+g|^p d\mu)^{1/p} \leq (\int |f|^p d\mu)^{1/p} + (\int |g|^p d\mu)^{1/p}$.

$f \rightarrow (\int |f|^p d\mu)^{1/p}$ would define a norm if $(\int |f|^p d\mu)^{1/p} = 0$ implied $f = 0$. Unfortunately, this is not so because if $f = 0$ a.e. but is nonzero on a set of measure zero, $(\int |f|^p d\mu)^{1/p} = 0$ and this is not allowed. However, all the other properties of a norm are available and so a little thing like a set of measure zero will not prevent the consideration of L^p as a normed vector space if two functions in L^p which differ only on a set of measure zero are considered the same. That is, an element of L^p is really an equivalence class of functions where two functions are equivalent if they are equal a.e. With this convention, here is a definition.

Definition 12.1.7 Let $f \in L^p(\Omega)$. Define $\|f\|_p \equiv \|f\|_{L^p} \equiv (\int |f|^p d\mu)^{1/p}$.

Then with this definition and using the convention that elements in L^p are considered to be the same if they differ only on a set of measure zero, $\|\cdot\|_p$ is a norm on $L^p(\Omega)$ because if $\|f\|_p = 0$ then $f = 0$ a.e. and so f is considered to be the zero function because it differs from 0 only on a set of measure zero.

The following is an important definition.

Definition 12.1.8 A complete normed linear space is called a Banach¹ space.

L^p is a Banach space. This is the next big theorem which says that these L^p spaces are always complete.

Theorem 12.1.9 The following holds for $L^p(\Omega, \mathcal{F}, \mu)$, $p \geq 1$. If $\{f_n\}$ is a Cauchy sequence in $L^p(\Omega)$, then there exists $f \in L^p(\Omega)$ and a subsequence which converges a.e. to $f \in L^p(\Omega)$, and $\|f_n - f\|_p \rightarrow 0$.

Proof: Let $\{f_n\}$ be a Cauchy sequence in $L^p(\Omega)$. This means that for every $\varepsilon > 0$ there exists N such that if $n, m \geq N$, then $\|f_n - f_m\|_p < \varepsilon$. Now select a subsequence as follows. Let n_1 be such that $\|f_n - f_m\|_p < 2^{-1}$ whenever $n, m \geq n_1$. Let n_2 be such that $n_2 > n_1$ and $\|f_n - f_m\|_p < 2^{-2}$ whenever $n, m \geq n_2$. If n_1, \dots, n_k have been chosen, let $n_{k+1} > n_k$ and whenever $n, m \geq n_{k+1}$, $\|f_n - f_m\|_p < 2^{-(k+1)}$. The subsequence just mentioned is $\{f_{n_k}\}$. Thus

$$\begin{aligned} & \mu \left(\left\{ \omega : |f_{n_k}(\omega) - f_{n_{k+1}}(\omega)|^p > \left(\frac{2}{3}\right)^k \right\} \right) \equiv \mu(E_k) \\ & \leq \left(\frac{3}{2}\right)^k \int_{E_k} |f_{n_k}(\omega) - f_{n_{k+1}}(\omega)|^p d\mu \end{aligned}$$

¹These spaces are named after Stefan Banach, 1892-1945. Banach spaces are the basic item of study in the subject of functional analysis and will be considered later in this book.

There is a recent biography of Banach, R. Kato, *The Life of Stefan Banach*, (A. Kostant and W. Woyczyński, translators and editors) Birkhauser, Boston (1996). More information on Banach can also be found in a recent short article written by Douglas Henderson who is in the department of chemistry and biochemistry at BYU.

Banach was born in Austria, worked in Poland and died in the Ukraine but never moved. This is because borders kept changing. There is a rumor that he died in a German concentration camp which is apparently not true. It seems he died after the war of lung cancer.

He was an interesting character. He hated taking examinations so much that he did not receive his undergraduate university degree. Nevertheless, he did become a professor of mathematics due to his important research. He and some friends would meet in a cafe called the Scottish cafe where they wrote on the marble table tops until Banach's wife supplied them with a notebook which became the "Scottish notebook" and was eventually published.

$$\leq \left(\frac{3}{2}\right)^k \|f_{n_k} - f_{n_{k+1}}\|_{L^p}^p < \left(\frac{3}{2}\right)^k 2^{-kp} \leq \frac{3^k}{4^k}$$

Hence $\sum_k \mu(E_k) < \infty$ and so, by the Borel Cantelli lemma, Lemma 9.2.5 on Page 243, there is a set of measure zero N such that if $\omega \notin N$, then $|f_{n_k}(\omega) - f_{n_{k+1}}(\omega)| \leq \left(\frac{2}{3}\right)^{k/p}$. Since

$$\sum_k \left(\frac{2}{3}\right)^{k/p} < \infty, \{\mathcal{X}_{N^c}(\omega) f_{n_k}(\omega)\}_{k=1}^\infty$$

is a Cauchy sequence for all ω . Let it converge to $f(\omega)$, a measurable function since it is a limit of measurable functions. By Fatou's lemma, and the Minkowski inequality, Corollary 12.1.5, $\|f - f_{n_k}\|_p = \left(\int |f - f_{n_k}|^p d\mu\right)^{1/p} \leq$

$$\begin{aligned} \liminf_{m \rightarrow \infty} \left(\int |f_{n_m} - f_{n_k}|^p d\mu \right)^{1/p} &= \liminf_{m \rightarrow \infty} \|f_{n_m} - f_{n_k}\|_p \leq \\ \liminf_{m \rightarrow \infty} \sum_{j=k}^{m-1} \|f_{n_{j+1}} - f_{n_j}\|_p &\leq \sum_{i=k}^\infty \|f_{n_{i+1}} - f_{n_i}\|_p \leq 2^{-(k-1)}. \end{aligned} \quad (12.3)$$

Therefore, $f \in L^p(\Omega)$ because $\|f\|_p \leq \|f - f_{n_k}\|_p + \|f_{n_k}\|_p < \infty$, and $\lim_{k \rightarrow \infty} \|f_{n_k} - f\|_p = 0$. This proves b.).

This has shown f_{n_k} converges to f in $L^p(\Omega)$. It follows the original Cauchy sequence also converges to f in $L^p(\Omega)$. This is a general fact that if a subsequence of a Cauchy sequence converges, then so does the original Cauchy sequence. This is Theorem 3.2.2. ■

In working with the L^p spaces, the following inequality also known as Minkowski's inequality is very useful. See [25]. It is similar to the Minkowski inequality for sums. To see this, replace the integral, \int_X with a finite summation sign and you will see the usual Minkowski inequality or rather the version of it given in Corollary 12.1.6.

Lemma 12.1.10 *Let (X, \mathcal{S}, μ) and $(Y, \mathcal{F}, \lambda)$ be finite measure spaces and let f be $\mu \times \lambda$ measurable. Then the following inequality is valid for $p \geq 1$.*

$$\int_X \left(\int_Y |f(x, y)|^p d\lambda \right)^{\frac{1}{p}} d\mu \geq \left(\int_Y \left(\int_X |f(x, y)| d\mu \right)^p d\lambda \right)^{\frac{1}{p}}. \quad (12.4)$$

Proof: This is an application of the Fubini theorem and Holder inequality. Recall that $p - 1 = p/p'$. Let $J(y) \equiv \int_X |f(x, y)| d\mu$. Then

$$\begin{aligned} \int_Y \left(\int_X |f(x, y)| d\mu \right)^p d\lambda &= \int_Y J(y)^{p/p'} \int_X |f(x, y)| d\mu d\lambda \\ &= \int_Y \int_X |f(x, y)| J(y)^{p/p'} d\mu d\lambda = \int_X \int_Y |f(x, y)| J(y)^{p/p'} d\lambda d\mu \\ &\leq \left(\int_Y J(y)^p d\lambda \right)^{1/p'} \int_X \left(\int_Y |f(x, y)|^p d\lambda \right)^{1/p} d\mu \\ &= \left(\int_Y \left(\int_X |f(x, y)| d\mu \right)^p d\lambda \right)^{1/p'} \int_X \left(\int_Y |f(x, y)|^p d\lambda \right)^{1/p} d\mu \end{aligned}$$

Thus $\int_Y (\int_X |f(x,y)| d\mu)^p d\lambda \leq$

$$\left(\int_Y \left(\int_X |f(x,y)| d\mu \right)^p d\lambda \right)^{1/p'} \int_X \left(\int_Y |f(x,y)|^p d\lambda \right)^{1/p} d\mu \quad (12.5)$$

If f is bounded, divide both sides by the first factor on the right and obtain 12.4. Otherwise replace f with $\min(f, n)$, divide and then apply the monotone convergence theorem as $n \rightarrow \infty$ to get 12.4. Note that 12.4 holds even if the first factor on the right in 12.5 equals zero. ■

Now consider the case where the measure spaces are σ finite.

Theorem 12.1.11 *Let (X, \mathcal{S}, μ) and $(Y, \mathcal{T}, \lambda)$ be σ -finite measure spaces and let f be product measurable. Then the following inequality is valid for $p \geq 1$.*

$$\int_X \left(\int_Y |f(x,y)|^p d\lambda \right)^{1/p} d\mu \geq \left(\int_Y \left(\int_X |f(x,y)| d\mu \right)^p d\lambda \right)^{1/p}. \quad (12.6)$$

Proof: Since the two measure spaces are σ finite, there exist measurable sets, X_m and Y_k such that $X_m \subseteq X_{m+1}$ for all m , $Y_k \subseteq Y_{k+1}$ for all k , and $\mu(X_m), \lambda(Y_k) < \infty$. From the above,

$$\int_{X_m} \left(\int_{Y_k} |f(x,y)|^p d\lambda \right)^{1/p} d\mu \geq \left(\int_{Y_k} \left(\int_{X_m} |f(x,y)| d\mu \right)^p d\lambda \right)^{1/p}. \quad (12.7)$$

Now use the monotone convergence theorem to pass to a limit first as $k \rightarrow \infty$ and then as $m \rightarrow \infty$. ■

Note that the proof of this theorem depends on two manipulations, the interchange of the order of integration and Holder's inequality. Also observe that there is nothing to check in the case of double sums. Thus if $a_{ij} \geq 0$, it is always the case that $(\sum_j (\sum_i a_{ij})^p)^{1/p} \leq \sum_i (\sum_j a_{ij}^p)^{1/p}$ because the integrals in this case are just sums and $(i, j) \rightarrow a_{ij}$ is measurable.

The L^p spaces have many important properties. Before considering these, here is a definition of L^∞

Definition 12.1.12 *$f \in L^\infty(\Omega, \mu)$ if there exists a set of measure zero E , and a constant $C < \infty$ such that $|f(x)| \leq C$ for all $x \notin E$. Then $\|f\|_\infty$ is defined as $\|f\|_\infty \equiv \inf\{C : |f(x)| \leq C \text{ a.e.}\}$, the inf of all such C .*

Proposition 12.1.13 *$\|\cdot\|_\infty$ is a norm on $L^\infty(\Omega, \mu)$ provided f and g are identified if $f(x) = g(x)$ a.e. Also, $L^\infty(\Omega, \mu)$ is complete. In addition to this, $\|f\|_\infty$ has the property that $|f(x)| \leq \|f\|_\infty$ for a.e. x so the norm is the smallest constant with this property.*

Proof: It is obvious that $\|f\|_\infty \geq 0$. Let $C_n \downarrow \|f\|_\infty$ where $|f(x)| \leq C_n$ off a set of measure zero E_n . Then let $E \equiv \cup_n E_n$. This is also a set of measure zero and if $x \notin E$, then $|f(x)| \leq C_n$ for all C_n and so $|f(x)| \leq \|f\|_\infty$ for all $x \notin E$. Thus $\|f\|_\infty$ is the smallest number C with $|f(x)| \leq C$ a.e. In case $\|f\|_\infty = 0$, $f(x) = 0$ off of E and so we regard f as 0 because it equals 0 a.e.

If $c = 0$ there is nothing to show in the claim that $\|cf\|_\infty = |c| \|f\|_\infty$. Assume then that $c \neq 0$. $|c| \|f\|_\infty \geq |cf(x)|$ a.e. Thus $|c| \|f\|_\infty \geq \|cf\|_\infty$ whenever $c \neq 0$. Thus $\|f\|_\infty = \|\frac{1}{c} cf\|_\infty \leq \frac{1}{|c|} \|cf\|_\infty$ and so $|c| \|f\|_\infty \leq \|cf\|_\infty$. It remains to verify the triangle inequality.

It was just shown that $\|f\|_\infty$ is the smallest constant such that $|f(x)| \leq \|f\|_\infty$ a.e. Hence if $f, g \in L^\infty(\Omega)$, $|f(x) + g(x)| \leq |f(x)| + |g(x)| \leq \|f\|_\infty + \|g\|_\infty$ a.e. and so by definition $\|f + g\|_\infty \leq \|f\|_\infty + \|g\|_\infty$.

Next suppose you have a Cauchy sequence in $L^\infty(\Omega)$ $\{f_n\}$. Let $|f_n(x) - f_m(x)| < \|f_n - f_m\|_\infty$ for $x \notin E_{nm}$, $\mu(E_{nm}) = 0$ and let $|f_n(x)| \leq \|f_n\|_\infty$ for $x \notin E_n$, $\mu(E_n) = 0$. Then let $E \equiv \bigcup_n E_n \cup \bigcup_{m,n} E_{nm}$. It follows that for $x \notin E$, $\lim_{n \rightarrow \infty} f_n(x)$ exists. Let $f(x)$ be this limit for $x \notin E$ and let $f(x) = 0$ on E . Also $|\|f_n\|_\infty - \|f_m\|_\infty| \leq \|f_n - f_m\|_\infty$ since $\|\cdot\|_\infty$ is a norm. Therefore, for $x \notin E$, $|f(x)| = \lim_{n \rightarrow \infty} |f_n(x)| \leq \lim_{n \rightarrow \infty} \|f_n\|_\infty \equiv C$ so $f \in L^\infty(\Omega, \mu)$. Also, for $x \notin E$, $|f_m(x) - f_n(x)| \leq \|f_m - f_n\|_\infty < \varepsilon$ if $m > n$ and n is large enough. Therefore, for such n , letting $m \rightarrow \infty$, $|f(x) - f_n(x)| \leq \varepsilon$ for $x \notin E$. It follows that $\|f - f_n\|_\infty \leq \varepsilon$ if n large enough and so by definition, $\lim_{n \rightarrow \infty} \|f - f_n\|_\infty = 0$. ■

12.2 Density Considerations

Theorem 12.2.1 *Let $p \geq 1$ and let $(\Omega, \mathcal{S}, \mu)$ be a measure space. Then the simple functions are dense in $L^p(\Omega)$. In fact, if $f \in L^p(\Omega)$, then there is a sequence of simple functions $\{s_n\}$ such that $|s_n| \leq |f|$ and $\|f - s_n\|_p \rightarrow 0$.*

Proof: Recall that a function f , having values in \mathbb{R} can be written in the form $f = f^+ - f^-$ where

$$f^+ = \max(0, f), \quad f^- = \max(0, -f).$$

Therefore, an arbitrary complex valued function, f is of the form

$$f = \operatorname{Re} f^+ - \operatorname{Re} f^- + i(\operatorname{Im} f^+ - \operatorname{Im} f^-).$$

If each of these nonnegative functions is approximated by a simple function, it follows f is also approximated by a simple function. Approximating each of the positive and negative parts with simple functions having absolute value less than what is approximated, it would follow that $|s_n| \leq 4|f|$ and all that is left is to verify that $\|s_n - f\|_p \rightarrow 0$ which occurs if it happens for each of these positive and negative parts of real and imaginary parts.

Now $|f(x) - s_n(x)| \leq 5|f|$ and so $|f(x) - s_n(x)|^p \leq 5^p |f|^p$ which is in L^1 . Then by the dominated convergence theorem, $0 = \lim_{n \rightarrow \infty} \int |f(x) - s_n(x)|^p d\mu$ showing that the simple functions are dense in L^p . ■

Note how this observation always holds and requires no assumptions on the measures.

Recall that for Ω a topological space, $C_c(\Omega)$ is the space of continuous functions with compact support in Ω . Also recall the following definition.

Definition 12.2.2 *Let $(\Omega, \mathcal{S}, \mu)$ be a measure space and suppose (Ω, τ) is also a topological space (metric space if you like.). Then $(\Omega, \mathcal{S}, \mu)$ is called a regular measure space if the σ algebra of Borel sets is contained in \mathcal{S} and for all $E \in \mathcal{S}$,*

$$\mu(E) = \inf\{\mu(V) : V \supseteq E \text{ and } V \text{ open}\}$$

and if $\mu(E) < \infty$,

$$\mu(E) = \sup\{\mu(K) : K \subseteq E \text{ and } K \text{ is compact}\}$$

and $\mu(K) < \infty$ for any compact set, K .

For example Lebesgue measure is an example of such a measure. More generally these measures are often referred to as Radon measures when they are complete. Recall the following important result which is Lemma 3.12.4.

Lemma 12.2.3 *Let Ω be a metric space in which the closed balls are compact and let K be a compact subset of V , an open set. Then there exists a continuous function $f : \Omega \rightarrow [0, 1]$ such that $f(x) = 1$ for all $x \in K$ and $\text{spt}(f)$ is a compact subset of V . That is, $K \prec f \prec V$.*

It is not necessary to be in a metric space to do this. You can accomplish the same thing using Urysohn's lemma in a normal topological space or, as is often done, a locally compact Hausdorff space. This can be discussed later.

Theorem 12.2.4 *Let $(\Omega, \mathcal{S}, \mu)$ be a regular measure space as in Definition 12.2.2 where the conclusion of Lemma 3.12.4 holds. Then $C_c(\Omega)$ is dense in $L^p(\Omega)$.*

Proof: First consider a measurable set, E where $\mu(E) < \infty$. Let $K \subseteq E \subseteq V$ where $\mu(V \setminus K) < \varepsilon$. Now let $K \prec h \prec V$. Then

$$\int |h - \mathcal{X}_E|^p d\mu \leq \int \mathcal{X}_{V \setminus K}^p d\mu = \mu(V \setminus K) < \varepsilon.$$

It follows that for each s a simple function in $L^p(\Omega)$, there exists $h \in C_c(\Omega)$ such that $\|s - h\|_p < \varepsilon$. This is because if $s(x) = \sum_{i=1}^m c_i \mathcal{X}_{E_i}(x)$ is a simple function in L^p where the c_i are the distinct nonzero values of s each $\mu(E_i) < \infty$ since otherwise $s \notin L^p$ due to the inequality $\int |s|^p d\mu \geq |c_i|^p \mu(E_i)$. By Theorem 12.2.1, simple functions are dense in $L^p(\Omega)$. Therefore, $C_c(\Omega)$ is dense in $L^p(\Omega)$. ■

12.3 Separability

The most important case is of course Lebesgue measure on \mathbb{R}^n or more generally, some Radon measure.

Theorem 12.3.1 *For $p \geq 1$ and μ a Radon measure, $L^p(\mathbb{R}^n, \mu)$ is separable. Recall this means there exists a countable set \mathcal{D} such that if $f \in L^p(\mathbb{R}^n, \mu)$ and $\varepsilon > 0$, there exists $g \in \mathcal{D}$ such that $\|f - g\|_p < \varepsilon$.*

Proof: Let Q be all functions of the form $c\mathcal{X}_{[a,b]}$ where

$$[a, b] \equiv [a_1, b_1] \times [a_2, b_2] \times \cdots \times [a_n, b_n],$$

and both a_i, b_i are rational, while c has rational real and imaginary parts. Let \mathcal{D} be the set of all finite sums of functions in Q . Thus, \mathcal{D} is countable. In fact \mathcal{D} is dense in $L^p(\mathbb{R}^n, \mu)$. To prove this, it is necessary to show that for every $f \in L^p(\mathbb{R}^n, \mu)$, there exists an element of \mathcal{D} , s such that $\|s - f\|_p < \varepsilon$. If it can be shown that for every $g \in C_c(\mathbb{R}^n)$ there exists $h \in \mathcal{D}$ such that $\|g - h\|_p < \varepsilon$, then this will suffice because if $f \in L^p(\mathbb{R}^n)$ is arbitrary, Theorem 12.2.4 implies there exists $g \in C_c(\mathbb{R}^n)$ such that $\|f - g\|_p \leq \frac{\varepsilon}{2}$ and then there would exist $h \in C_c(\mathbb{R}^n)$ such that $\|h - g\|_p < \frac{\varepsilon}{2}$. By the triangle inequality,

$$\|f - h\|_p \leq \|h - g\|_p + \|g - f\|_p < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon.$$

Therefore, assume at the outset that $f \in C_c(\mathbb{R}^n)$.

Let \mathcal{P}_m consist of all sets of the form $[\mathbf{a}, \mathbf{b}) \equiv \prod_{i=1}^n [a_i, b_i)$ where $a_i = j2^{-m}$ and $b_i = (j+1)2^{-m}$ for j an integer. Thus \mathcal{P}_m consists of a tiling of \mathbb{R}^n into half open rectangles having diameters $2^{-m}n^{\frac{1}{2}}$. There are countably many of these rectangles; so let $\mathcal{P}_m = \{[\mathbf{a}_i, \mathbf{b}_i)\}$ for $i \geq 1$, and $\mathbb{R}^n = \cup_{i=1}^\infty [\mathbf{a}_i, \mathbf{b}_i)$. Let c_i^m be complex numbers with rational real and imaginary parts satisfying

$$|f(\mathbf{a}_i) - c_i^m| < 2^{-m}, \quad |c_i^m| \leq |f(\mathbf{a}_i)|. \quad (12.8)$$

Let $s_m(\mathbf{x}) = \sum_{i=1}^\infty c_i^m \chi_{[\mathbf{a}_i, \mathbf{b}_i)}(\mathbf{x})$.

Since $f(\mathbf{a}_i) = 0$ except for finitely many values of i , the above is a finite sum. Then 12.8 implies $s_m \in \mathcal{D}$. If s_m converges uniformly to f then it will follow that s_m is close to f in L^p .

Since $f \in C_c(\mathbb{R}^n)$ it follows that f is uniformly continuous and so given $\varepsilon > 0$ there exists $\delta > 0$ such that if $|\mathbf{x} - \mathbf{y}| < \delta$, $|f(\mathbf{x}) - f(\mathbf{y})| < \varepsilon/2$. Now let m be large enough that every box in \mathcal{P}_m has diameter less than δ and also that $2^{-m} < \varepsilon/2$. Then if $[\mathbf{a}_i, \mathbf{b}_i)$ is one of these boxes of \mathcal{P}_m , and $\mathbf{x} \in [\mathbf{a}_i, \mathbf{b}_i)$, $|f(\mathbf{x}) - f(\mathbf{a}_i)| < \varepsilon/2$ and $|f(\mathbf{a}_i) - c_i^m| < 2^{-m} < \varepsilon/2$. Therefore, using the triangle inequality, it follows that for $\mathbf{x} \in [\mathbf{a}_i, \mathbf{b}_i)$,

$$\begin{aligned} |f(\mathbf{x}) - s_m(\mathbf{x})| &= |f(\mathbf{x}) - c_i^m| = |f(\mathbf{x}) - f(\mathbf{a}_i)| + |f(\mathbf{a}_i) - c_i^m| \\ &< 2^{-m} + 2^{-m} < \varepsilon \end{aligned}$$

and since \mathbf{x} is arbitrary, this establishes uniform convergence. From the construction, s_m and f are zero off some compact set K which does not depend on m . Therefore, for m large, $\int |f(\mathbf{x}) - s_m(\mathbf{x})|^p d\mu < \varepsilon^p \mu(K)$ and since ε is arbitrary, this shows that the countable set \mathcal{D} is dense in $L^p(\Omega)$ as claimed. ■

Here is an easier proof if you know the Weierstrass approximation theorem, Theorem 5.7.1 for example.

Theorem 12.3.2 For $p \geq 1$ and μ a Radon measure, $L^p(\mathbb{R}^n, \mu)$ is separable. Recall this means there exists a countable set \mathcal{D} , such that if $f \in L^p(\mathbb{R}^n, \mu)$ and $\varepsilon > 0$, there exists $g \in \mathcal{D}$ such that $\|f - g\|_p < \varepsilon$.

Proof: As noted above, the continuous functions with compact support are dense in $L^p(\mathbb{R}^n, \mu)$. Let $\text{spt}(f) \subseteq (-R, R)^n$. Consider the polynomials having rational coefficients. By the Weierstrass approximation theorem, and adjusting coefficients to make them all rational, there exists p , $\|f - p\|_{[-R, R]^n} < \frac{\varepsilon}{R^{n/p}}$ the norm being the uniform norm $\|g\|_{[-R, R]^n} \equiv \max \{|g(\mathbf{x})| : \mathbf{x} \in [-R, R]^n\}$. Now let $\text{spt}(f) \prec \tau_\varepsilon \prec V_\varepsilon$ where $\mu(V \setminus \text{spt}(f)) < \varepsilon^p$. Then

$$\begin{aligned} \|f - \tau_\varepsilon p\|_p^p &\leq \int_{\text{spt}(f)} \left(\frac{\varepsilon}{R^{n/p}}\right)^p d\mu + \int_{V \setminus \text{spt}(f)} d\mu \\ &\leq \frac{\varepsilon^p}{R^n} 4^n R^n + \mu(V \setminus \text{spt}(f)) < (1 + 4^n) \varepsilon^p \end{aligned}$$

Letting \mathcal{P} denote the polynomials with rational coefficients, let $\varepsilon_k \rightarrow 0$ and consider the set of functions $\cup_k \tau_{\varepsilon_k} \mathcal{P} \equiv \mathcal{D}$. This is countable and from the above computation, it is dense in L^p . ■

Corollary 12.3.3 Let Ω be any μ measurable subset of \mathbb{R}^n and let μ be a Radon measure. Then $L^p(\Omega, \mu)$ is separable. Here the σ algebra of measurable sets will consist of all intersections of measurable sets with Ω and the measure will be μ restricted to these sets.

Proof: Let $\tilde{\mathcal{D}}$ be the restrictions of \mathcal{D} to Ω . If $f \in L^p(\Omega)$, let F be the zero extension of f to all of \mathbb{R}^n . Let $\varepsilon > 0$ be given. By Theorem 12.3.1 or 12.3.2 there exists $s \in \mathcal{D}$ such that $\|F - s\|_p < \varepsilon$. Thus

$$\|s - f\|_{L^p(\Omega, \mu)} \leq \|s - F\|_{L^p(\mathbb{R}^n, \mu)} < \varepsilon$$

and so the countable set $\tilde{\mathcal{D}}$ is dense in $L^p(\Omega)$. ■

12.4 Continuity of Translation

Definition 12.4.1 Let f be a function defined on $U \subseteq \mathbb{R}^n$ and let $\mathbf{w} \in \mathbb{R}^n$. Then $f_{\mathbf{w}}$ will be the function defined on $\mathbf{w} + U$ by $f_{\mathbf{w}}(\mathbf{x}) = f(\mathbf{x} - \mathbf{w})$.

Theorem 12.4.2 (Continuity of translation in L^p) Let $f \in L^p(\mathbb{R}^n)$ with the measure being Lebesgue measure. Then $\lim_{\|\mathbf{w}\| \rightarrow 0} \|f_{\mathbf{w}} - f\|_p = 0$.

Proof: Let $\varepsilon > 0$ be given and let $g \in C_c(\mathbb{R}^n)$ with $\|g - f\|_p < \frac{\varepsilon}{3}$. Since Lebesgue measure is translation invariant ($m_n(\mathbf{w} + E) = m_n(E)$), $\|g_{\mathbf{w}} - f_{\mathbf{w}}\|_p = \|g - f\|_p < \frac{\varepsilon}{3}$. You can see this from looking at simple functions and passing to the limit or you could use the change of variables formula to verify it.

Therefore

$$\|f - f_{\mathbf{w}}\|_p \leq \|f - g\|_p + \|g - g_{\mathbf{w}}\|_p + \|g_{\mathbf{w}} - f_{\mathbf{w}}\|_p < \frac{2\varepsilon}{3} + \|g - g_{\mathbf{w}}\|_p. \quad (12.9)$$

But $\lim_{|\mathbf{w}| \rightarrow 0} g_{\mathbf{w}}(\mathbf{x}) = g(\mathbf{x})$ uniformly in \mathbf{x} because g is uniformly continuous. Now let B be a large ball containing $\text{spt}(g)$ and let δ_1 be small enough that $B(\mathbf{x}, \delta) \subseteq B$ whenever $\mathbf{x} \in \text{spt}(g)$. If $\varepsilon > 0$ is given there exists $\delta < \delta_1$ such that if $|\mathbf{w}| < \delta$, it follows that $|g(\mathbf{x} - \mathbf{w}) - g(\mathbf{x})| < \varepsilon/3 \left(1 + m_n(B)^{1/p}\right)$. Therefore,

$$\|g - g_{\mathbf{w}}\|_p = \left(\int_B |g(\mathbf{x}) - g(\mathbf{x} - \mathbf{w})|^p dm_n \right)^{1/p} \leq \varepsilon \frac{m_n(B)^{1/p}}{3(1 + m_n(B)^{1/p})} < \frac{\varepsilon}{3}.$$

Thus, whenever $|\mathbf{w}| < \delta$, it follows $\|g - g_{\mathbf{w}}\|_p < \frac{\varepsilon}{3}$ and so from 12.9 $\|f - f_{\mathbf{w}}\|_p < \varepsilon$. ■

Here is a remarkable corollary.

Corollary 12.4.3 Suppose $f \in L^1(\mathbb{R}^p, m_p)$ and let \mathbf{v} be any nonzero vector. Then there is a set of measure zero N and a sequence $t_n \rightarrow 0+$ such that if $\mathbf{x} \notin N$, and $0 < s_n \leq t_n$

$$\lim_{n \rightarrow \infty} |f(\mathbf{x}) - f(\mathbf{x} + s_n \mathbf{v})| = 0.$$

Proof: Let t_n be such that if $s_n \leq t_n$, $\|f - f(\cdot + s_n \mathbf{v})\|_{L^1} < 4^{-n}$. This exists by continuity of translation in $L^1(\mathbb{R}^p, m_p)$. Then

$$m_p(E_n) \equiv m_p(\{\mathbf{x} : |f(\mathbf{x}) - f(\mathbf{x} + s_n \mathbf{v})| \geq 2^{-n}\}) \leq \frac{\int |f - f(\cdot + s_n \mathbf{v})| dm_p}{2^{-n}} < 2^{-n}$$

Thus there is a set of measure zero N such that if $\mathbf{x} \notin N$, then \mathbf{x} is in only finitely many of the sets E_n . It follows that for all n sufficiently large $|f(\mathbf{x}) - f(\mathbf{x} + s_n \mathbf{v})| < 2^{-n}$. ■

12.5 Mollifiers and Density of Smooth Functions

Definition 12.5.1 Let U be an open subset of \mathbb{R}^n . $C_c^\infty(U)$ is the vector space of all infinitely differentiable functions which equal zero for all x outside of some compact set contained in U . Similarly, $C_c^m(U)$ is the vector space of all functions which are m times continuously differentiable and whose support is a compact subset of U .

Example 12.5.2 Let $U = B(z, 2r)$

$$\psi(x) = \begin{cases} \exp \left[\left(|x - z|^2 - r^2 \right)^{-1} \right] & \text{if } |x - z| < r, \\ 0 & \text{if } |x - z| \geq r. \end{cases}$$

Then a little work shows $\psi \in C_c^\infty(U)$. The following also is easily obtained.

Lemma 12.5.3 Let U be any open set. Then $C_c^\infty(U) \neq \emptyset$.

Proof: Pick $z \in U$ and let r be small enough that $B(z, 2r) \subseteq U$. Then let

$$\psi \in C_c^\infty(B(z, 2r)) \subseteq C_c^\infty(U)$$

be the function of the above example.

Definition 12.5.4 Let $U = \{x \in \mathbb{R}^n : |x| < 1\}$. A sequence $\{\psi_m\} \subseteq C_c^\infty(U)$ is called a mollifier² if $\psi_m(x) \geq 0$, $\psi_m(x) = 0$, if $|x| \leq \frac{1}{m}$, and $\int \psi_m(x) = 1$. Sometimes it may be written as $\{\psi_\varepsilon\}$ where ψ_ε satisfies the above conditions except $\psi_\varepsilon(x) = 0$ if $|x| \geq \varepsilon$. In other words, ε takes the place of $1/m$ and in everything that follows $\varepsilon \rightarrow 0$ instead of $m \rightarrow \infty$.

As before, $\int f(x, y) d\mu(y)$ will mean x is fixed and the function $y \rightarrow f(x, y)$ is being integrated. To make the notation more familiar, dx is written instead of $dm_n(x)$.

Example 12.5.5 Let $\psi \in C_c^\infty(B(0, 1))$ with $\psi(x) \geq 0$ and $\int \psi dm = 1$. Let $\psi_m(x) = c_m \psi(mx)$ where c_m is chosen in such a way that $\int \psi_m dm = 1$. By the change of variables theorem $c_m = m^n$. Also ψ_m is zero off $B(0, 1/m)$.

Definition 12.5.6 A function f , is said to be in $L_{loc}^1(\mathbb{R}^n, \mu)$ if f is μ measurable and if $|f| \chi_K \in L^1(\mathbb{R}^n, \mu)$ for every compact set K . Here μ is a regular, complete measure on \mathbb{R}^n . Usually $\mu = m_n$, Lebesgue measure. When this is so, write $L_{loc}^1(\mathbb{R}^n)$, etc. If $f \in L_{loc}^1(\mathbb{R}^n, \mu)$, and $g \in C_c(\mathbb{R}^n)$, $f * g(x) \equiv \int f(y)g(x - y)d\mu$.

The following lemma will be useful in what follows. It says that one of these very un-regular functions in $L_{loc}^1(\mathbb{R}^n, \mu)$ is smoothed out by convolving with a mollifier.

Lemma 12.5.7 Let $f \in L_{loc}^1(\mathbb{R}^n, \mu)$, and $g \in C_c^\infty(\mathbb{R}^n)$. Then $f * g$ is an infinitely differentiable function. Here μ is a Radon measure on \mathbb{R}^n . In case f is continuous with compact support $\text{spt}(f)$, and if ψ_m is a mollifier as described above, then $\text{spt}(f * \psi_m) \subseteq \text{spt}(f) + B(0, 1/m)$. Also $\|f - f * \psi_m\| \rightarrow 0$.

²This is sometimes called an approximate identity if the differentiability is not included.

Proof: Consider the difference quotient for calculating a partial derivative of $f * g$.

$$\frac{f * g(\mathbf{x} + t\mathbf{e}_j) - f * g(\mathbf{x})}{t} = \int f(\mathbf{y}) \frac{g(\mathbf{x} + t\mathbf{e}_j - \mathbf{y}) - g(\mathbf{x} - \mathbf{y})}{t} d\mu(\mathbf{y}).$$

Using the fact that $g \in C_c^\infty(\mathbb{R}^n)$, the quotient $\frac{g(\mathbf{x} + t\mathbf{e}_j - \mathbf{y}) - g(\mathbf{x} - \mathbf{y})}{t}$ is uniformly bounded. To see this easily, use Theorem 7.5.2 on Page 190 to get the existence of a constant, M depending on $\max\{|Dg(\mathbf{x})| : \mathbf{x} \in \mathbb{R}^n\}$ such that $|g(\mathbf{x} + t\mathbf{e}_j - \mathbf{y}) - g(\mathbf{x} - \mathbf{y})| \leq M|t|$ for any choice of \mathbf{x} and \mathbf{y} . Therefore, there exists a dominating function for the integrand of the above integral which is of the form $C|f(\mathbf{y})|\mathcal{X}_K$ where K is a compact set depending on the support of g . It follows the limit of the difference quotient above passes inside the integral as $t \rightarrow 0$ and $\frac{\partial}{\partial x_j}(f * g)(\mathbf{x}) = \int f(\mathbf{y}) \frac{\partial}{\partial x_j} g(\mathbf{x} - \mathbf{y}) d\mu(\mathbf{y})$. Now letting $\frac{\partial}{\partial x_j} g$ play the role of g in the above argument, partial derivatives of all orders exist. A similar use of the dominated convergence theorem shows all these partial derivatives are also continuous.

For the last claim, it is clear that $\text{spt}(f * \psi_m) \subseteq \text{spt}(f) + B(\mathbf{0}, 1/m)$ since off $\text{spt}(f) + B(\mathbf{0}, 1/m)$ the integral for $f * \psi_m$ will be 0. To verify the last claim, let $\varepsilon > 0$ be given. By uniform continuity of f , $|f(\mathbf{x}) - f(\mathbf{x} - \mathbf{y})| < \varepsilon$ whenever $|\mathbf{y}|$ is sufficiently small. Therefore,

$$\begin{aligned} |f(\mathbf{x}) - f * \psi_m(\mathbf{x})| &= \left| \int (f(\mathbf{x}) - f(\mathbf{x} - \mathbf{y})) \psi_m(\mathbf{y}) d\mu(\mathbf{y}) \right| \\ &\leq \int_{B(\mathbf{0}, 1/m)} |f(\mathbf{x}) - f(\mathbf{x} - \mathbf{y})| \psi_m(\mathbf{y}) d\mu(\mathbf{y}) < \varepsilon \int \psi_m d\mu = \varepsilon \end{aligned}$$

whenever m is large enough. ■

Theorem 12.5.8 For each $p \geq 1$, $C_c^\infty(\mathbb{R}^n)$ is dense in $L^p(\mathbb{R}^n)$. Here the measure is Lebesgue measure.

Proof: Let $f \in L^p(\mathbb{R}^n)$ and let $\varepsilon > 0$ be given. Choose $g \in C_c(\mathbb{R}^n)$ such that $\|f - g\|_p < \frac{\varepsilon}{2}$. This can be done by using Theorem 12.2.4. Now let

$$g_m(\mathbf{x}) = g * \psi_m(\mathbf{x}) \equiv \int g(\mathbf{x} - \mathbf{y}) \psi_m(\mathbf{y}) dm_n(\mathbf{y}) = \int g(\mathbf{y}) \psi_m(\mathbf{x} - \mathbf{y}) dm_n(\mathbf{y})$$

where $\{\psi_m\}$ is a mollifier. It follows from Lemma 12.5.7 $g_m \in C_c^\infty(\mathbb{R}^n)$. It vanishes if $\mathbf{x} \notin \text{spt}(g) + B(\mathbf{0}, \frac{1}{m})$.

$$\begin{aligned} \|g - g_m\|_p &= \left(\int |g(\mathbf{x}) - \int g(\mathbf{x} - \mathbf{y}) \psi_m(\mathbf{y}) dm_n(\mathbf{y})|^p dm_n(\mathbf{x}) \right)^{\frac{1}{p}} \\ &\leq \left(\int \left(\int |g(\mathbf{x}) - g(\mathbf{x} - \mathbf{y})| \psi_m(\mathbf{y}) dm_n(\mathbf{y}) \right)^p dm_n(\mathbf{x}) \right)^{\frac{1}{p}} \\ &\leq \int \left(\int |g(\mathbf{x}) - g(\mathbf{x} - \mathbf{y})|^p dm_n(\mathbf{x}) \right)^{\frac{1}{p}} \psi_m(\mathbf{y}) dm_n(\mathbf{y}) \\ &= \int_{B(\mathbf{0}, \frac{1}{m})} \|g - g_{\mathbf{y}}\|_p \psi_m(\mathbf{y}) dm_n(\mathbf{y}) < \frac{\varepsilon}{2} \end{aligned}$$

whenever m is large enough thanks to the uniform continuity of g . Theorem 12.1.11 was used to obtain the third inequality. There is no measurability problem because the function

$$(\mathbf{x}, \mathbf{y}) \rightarrow |g(\mathbf{x}) - g(\mathbf{x} - \mathbf{y})| \psi_m(\mathbf{y})$$

is continuous. Thus when m is large enough,

$$\|f - g_m\|_p \leq \|f - g\|_p + \|g - g_m\|_p < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon. \blacksquare$$

This is a very remarkable result. Functions in $L^p(\mathbb{R}^n)$ don't need to be continuous anywhere and yet every such function is very close in the L^p norm to one which is infinitely differentiable having compact support. The same result holds for $L^p(U)$ for U an open set. This is the next corollary.

Corollary 12.5.9 *Let U be an open set. For each $p \geq 1$, $C_c^\infty(U)$ is dense in $L^p(U)$. Here the measure is Lebesgue measure.*

Proof: Let $f \in L^p(U)$ and let $\varepsilon > 0$ be given. Choose $g \in C_c(U)$ such that $\|f - g\|_p < \frac{\varepsilon}{2}$. This is possible because Lebesgue measure restricted to the open set, U is regular. Thus the existence of such a g follows from Theorem 12.2.4. Now let

$$g_m(x) = g * \psi_m(x) \equiv \int g(x - y) \psi_m(y) dm_n(y) = \int g(y) \psi_m(x - y) dm_n(y)$$

where $\{\psi_m\}$ is a mollifier. It follows from Lemma 12.5.7 $g_m \in C_c^\infty(U)$ for all m sufficiently large. It vanishes if $x \notin \text{spt}(g) + B(0, \frac{1}{m})$. Then

$$\begin{aligned} \|g - g_m\|_p &= \left(\int |g(x) - \int g(x - y) \psi_m(y) dm_n(y)|^p dm_n(x) \right)^{\frac{1}{p}} \\ &\leq \left(\int \left(\int |g(x) - g(x - y)| \psi_m(y) dm_n(y) \right)^p dm_n(x) \right)^{\frac{1}{p}} \\ &\leq \int \left(\int |g(x) - g(x - y)|^p dm_n(x) \right)^{\frac{1}{p}} \psi_m(y) dm_n(y) \\ &= \int_{B(0, \frac{1}{m})} \|g - g_y\|_p \psi_m(y) dm_n(y) < \frac{\varepsilon}{2} \end{aligned}$$

whenever m is large enough thanks to uniform continuity of g . Theorem 12.1.11 was used to obtain the third inequality. There is no measurability problem because the function

$$(x, y) \rightarrow |g(x) - g(x - y)| \psi_m(y)$$

is continuous. Thus when m is large enough,

$$\|f - g_m\|_p \leq \|f - g\|_p + \|g - g_m\|_p < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon. \blacksquare$$

Another thing should probably be mentioned. If you have had a course in complex analysis, you may be wondering whether these infinitely differentiable functions having compact support have anything to do with analytic functions which also have infinitely many derivatives. The answer is no! Recall that if an analytic function has a limit point in the set of zeros then it is identically equal to zero. Thus these functions in $C_c^\infty(\mathbb{R}^n)$ are not analytic. This is a strictly real analysis phenomenon and has absolutely nothing to do with the theory of functions of a complex variable.

12.6 Smooth Partitions of Unity

Partitions of unity were discussed earlier. Here the idea of a smooth partition of unity is considered. The earlier general result on metric space is Theorem 3.12.5 on Page 92. Recall the following notation.

Notation 12.6.1 I will write $\phi \prec V$ to symbolize $\phi \in C_c(V)$, ϕ has values in $[0, 1]$, and ϕ has compact support in V . I will write $K \prec \phi \prec V$ for K compact and V open to symbolize ϕ is 1 on K and ϕ has values in $[0, 1]$ with compact support contained in V .

Definition 12.6.2 A collection of sets \mathcal{H} is called locally finite if for every x , there exists $r > 0$ such that $B(x, r)$ has nonempty intersection with only finitely many sets of \mathcal{H} . Of course every finite collection of sets is locally finite. This is the case of most interest in this book but the more general notion is interesting.

The thing about locally finite collection of sets is that the closure of their union equals the union of their closures. This is clearly true of a finite collection.

Lemma 12.6.3 Let \mathcal{H} be a locally finite collection of sets of a normed vector space V . Then

$$\overline{\cup \mathcal{H}} = \cup \{ \overline{H} : H \in \mathcal{H} \}.$$

Proof: It is obvious \supseteq holds in the above claim. It remains to go the other way. Suppose then that p is a limit point of $\cup \mathcal{H}$ and $p \notin \cup \mathcal{H}$. There exists $r > 0$ such that $B(p, r)$ has nonempty intersection with only finitely many sets of \mathcal{H} say these are H_1, \dots, H_m . Then I claim p must be a limit point of one of these. If this is not so, there would exist $r' > 0$ such that $0 < r' < r$ with $B(p, r')$ having empty intersection with each of these H_i . But then p would fail to be a limit point of $\cup \mathcal{H}$. Therefore, p is contained in the right side. It is clear $\cup \mathcal{H}$ is contained in the right side and so This proves the lemma. ■

A good example to consider is the rational numbers each being a set in \mathbb{R} . This is **not** a locally finite collection of sets and you note that $\overline{\mathbb{Q}} = \mathbb{R} \neq \cup \{ \overline{x} : x \in \mathbb{Q} \}$. By contrast, \mathbb{Z} is a locally finite collection of sets, the sets consisting of individual integers. The closure of \mathbb{Z} is equal to \mathbb{Z} because \mathbb{Z} has no limit points so it contains them all.

Lemma 12.6.4 Let K be a closed set in \mathbb{R}^p and let $\{V_i\}_{i=1}^\infty$ be a locally finite sequence of **bounded** open sets whose union contains K . Then there exist functions, $\psi_i \in C_c^\infty(V_i)$ such that for all $x \in K$, $1 = \sum_{i=1}^\infty \psi_i(x)$ and the function $f(x)$ given by $f(x) = \sum_{i=1}^\infty \psi_i(x)$ is in $C^\infty(\mathbb{R}^p)$.

Proof: Let $K_1 = K \setminus \cup_{i=2}^\infty V_i$. Thus K_1 is compact because it is a closed subset of a bounded set and $K_1 \subseteq V_1$. Let W_1 be an open set having compact closure which satisfies

$$K_1 \subseteq W_1 \subseteq \overline{W_1} \subseteq V_1$$

Thus W_1, V_2, \dots covers K and $\overline{W_1} \subseteq V_1$. Suppose W_1, \dots, W_r have been defined such that $\overline{W_i} \subseteq V_i$ for each i , and $W_1, \dots, W_r, V_{r+1}, \dots$ covers K . Then let

$$K_{r+1} \equiv K \setminus ((\cup_{i=r+2}^\infty V_i) \cup (\cup_{j=1}^r W_j)).$$

It follows K_{r+1} is compact because $K_{r+1} \subseteq V_{r+1}$. Let W_{r+1} satisfy

$$K_{r+1} \subseteq W_{r+1} \subseteq \overline{W_{r+1}} \subseteq V_{r+1}, \overline{W_{r+1}} \text{ is compact}$$

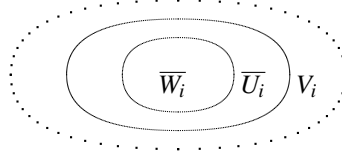
Continuing this way defines a sequence of open sets $\{W_i\}_{i=1}^\infty$ having compact closures with the property

$$\overline{W_i} \subseteq V_i, K \subseteq \bigcup_{i=1}^\infty W_i.$$

Note $\{W_i\}_{i=1}^\infty$ is locally finite because the original sequence, $\{V_i\}_{i=1}^\infty$ was locally finite. Now let U_i be open sets which satisfy

$$\overline{W_i} \subseteq U_i \subseteq \overline{U_i} \subseteq V_i, \overline{U_i} \text{ is compact.}$$

Similarly, $\{U_i\}_{i=1}^\infty$ is locally finite.



Now the local finiteness implies $\bigcup_{i=1}^\infty \overline{W_i} = \bigcup_{i=1}^\infty \overline{W_i}$. Define ϕ_i and γ , continuous having compact support such that

$$\overline{U_i} \prec \phi_i \prec V_i, \bigcup_{i=1}^\infty \overline{W_i} \prec \gamma \prec \bigcup_{i=1}^\infty U_i.$$

by convolving each of these with a mollifier, we can use Lemma 12.5.7 to preserve the above and also have each of these functions infinitely differentiable. Now define

$$\psi_i(x) = \begin{cases} \gamma(x)\phi_i(x)/\sum_{j=1}^\infty \phi_j(x) & \text{if } \sum_{j=1}^\infty \phi_j(x) \neq 0, \\ 0 & \text{if } \sum_{j=1}^\infty \phi_j(x) = 0. \end{cases}$$

All of these infinite sums are really finite sums because of the local finiteness of the $\{V_i\}$. Thus for \mathbf{y} near a given \mathbf{x} , all $\phi_j(\mathbf{y})$ are zero. Therefore, all continuity and differentiability of the individual ϕ_j is retained by the “infinite” sum.

If \mathbf{x} is such that $\sum_{j=1}^\infty \phi_j(\mathbf{x}) = 0$, then $\mathbf{x} \notin \bigcup_{i=1}^\infty \overline{U_i}$ because ϕ_i equals one on $\overline{U_i}$. Consequently $\gamma(\mathbf{y}) = 0$ for all \mathbf{y} near \mathbf{x} thanks to the fact that $\bigcup_{i=1}^\infty \overline{U_i}$ is closed and so $\psi_i(\mathbf{y}) = 0$ for all \mathbf{y} near \mathbf{x} . Hence ψ_i is infinitely differentiable at such \mathbf{x} . If $\sum_{j=1}^\infty \phi_j(\mathbf{x}) \neq 0$, this situation persists near \mathbf{x} because each ϕ_j is continuous and so ψ_i is infinitely differentiable at such points also. Therefore ψ_i is infinitely differentiable. If $\mathbf{x} \in K$, then $\gamma(\mathbf{x}) = 1$ and so $\sum_{j=1}^\infty \psi_j(\mathbf{x}) = 1$. Clearly $0 \leq \psi_i(\mathbf{x}) \leq 1$ and $\text{spt}(\psi_j) \subseteq V_j$. ■

The functions, $\{\psi_i\}$ are called a C^∞ partition of unity. The following is very useful.

Corollary 12.6.5 *In the context of Lemma 12.6.4, if H is a compact subset of V_i for some V_i there exists a partition of unity such that $\psi_i(x) = 1$ for all $x \in H$ in addition to the conclusion of Lemma 12.6.4.*

Proof: Keep V_i the same but replace all the V_j with $\tilde{V}_j \equiv V_j \setminus H$. Now in the proof above, applied to this modified collection of open sets, if $j \neq i$, $\phi_j(x) = 0$ whenever $x \in H$. Therefore, $\psi_i(x) = 1$ on H . ■

If K is compact, we can always reduce to a finite cover and so we obtain the following:

Theorem 12.6.6 *Let K be a compact set in \mathbb{R}^n and let $\{U_i\}_{i=1}^\infty$ be an open cover of K . Then there exist functions, $\psi_k \in C_c^\infty(U_i)$ such that $\psi_i \prec U_i$ and for all $\mathbf{x} \in K$, it follows that $\sum_{i=1}^\infty \psi_i(\mathbf{x}) = 1$. If K_1 is a compact subset of U_1 there exist such functions such that also $\psi_1(\mathbf{x}) = 1$ for all $\mathbf{x} \in K_1$.*

12.7 Exercises

1. Let E be a Lebesgue measurable set in \mathbb{R} . Suppose $m(E) > 0$. Consider the set $E - E = \{x - y : x \in E, y \in E\}$. Show that $E - E$ contains an interval. **Hint:** Let $f(x) = \int \mathcal{X}_E(t) \mathcal{X}_E(x+t) dt$. Note f is continuous at 0 and $f(0) > 0$ and use continuity of translation in L^p .
2. Establish the inequality $\|fg\|_r \leq \|f\|_p \|g\|_q$ whenever $\frac{1}{r} = \frac{1}{p} + \frac{1}{q}$.
3. Let $(\Omega, \mathcal{S}, \mu)$ be counting measure on \mathbb{N} . Thus $\Omega = \mathbb{N}$ and $\mathcal{S} = \mathcal{P}(\mathbb{N})$ with $\mu(S) =$ number of things in S . Let $1 \leq p \leq q$. Show that in this case, $L^1(\mathbb{N}) \subseteq L^p(\mathbb{N}) \subseteq L^q(\mathbb{N})$. **Hint:** This is real easy if you consider what $\int_{\Omega} f d\mu$ equals. How are the norms related?
4. Consider the function, $f(x, y) = \frac{x^{p-1}}{py} + \frac{y^{q-1}}{qx}$ for $x, y > 0$ and $\frac{1}{p} + \frac{1}{q} = 1$. Show directly that $f(x, y) \geq 1$ for all such x, y and show this implies $xy \leq \frac{x^p}{p} + \frac{y^q}{q}$.
5. Give an example of a sequence of functions in $L^p(\mathbb{R})$ which converges to zero in L^p but does not converge pointwise to 0. Does this contradict the proof of the theorem that L^p is complete?
6. Let K be a bounded subset of $L^p(\mathbb{R}^n)$ and suppose that there exists G such that \overline{G} is compact with $\int_{\mathbb{R}^n \setminus \overline{G}} |u(x)|^p dx < \varepsilon^p$ and for all $\varepsilon > 0$, there exist a $\delta > 0$ and such that if $|h| < \delta$, then $\int |u(x+h) - u(x)|^p dx < \varepsilon^p$ for all $u \in K$. Show that K is precompact in $L^p(\mathbb{R}^n)$. **Hint:** Let ϕ_k be a mollifier and consider $K_k \equiv \{u * \phi_k : u \in K\}$. Verify the conditions of the Ascoli Arzela theorem for these functions defined on \overline{G} and show there is an ε net for each $\varepsilon > 0$. Can you modify this to let an arbitrary open set take the place of \mathbb{R}^n ?
7. Let (Ω, d) be a metric space and suppose also that $(\Omega, \mathcal{S}, \mu)$ is a regular measure space such that $\mu(\Omega) < \infty$ and let $f \in L^1(\Omega)$ where f has complex values. Show that for every $\varepsilon > 0$, there exists an open set of measure less than ε , denoted here by V and a continuous function, g defined on Ω such that $f = g$ on V^C . Thus, aside from a set of small measure, f is continuous. If $|f(\omega)| \leq M$, show that it can be assumed that $|g(\omega)| \leq M$. This is called Lusin's theorem. **Hint:** Use Theorems 12.2.4 and 12.1.9 to obtain a sequence of functions in $C_c(\Omega)$, $\{g_n\}$ which converges pointwise a.e. to f and then use Egoroff's theorem to obtain a small set, W of measure less than $\varepsilon/2$ such that convergence is uniform on W^C . Now let F be a closed subset of W^C such that $\mu(W^C \setminus F) < \varepsilon/2$. Let $V = F^C$. Thus $\mu(V) < \varepsilon$ and on $F = V^C$, the convergence of $\{g_n\}$ is uniform showing that the restriction of f to V^C is continuous. Now use the Tietze extension theorem.
8. Let $\phi_m \in C_c^\infty(\mathbb{R}^n)$, $\phi_m(x) \geq 0$ and $\int_{\mathbb{R}^n} \phi_m(y) dy = 1$ with

$$\lim_{m \rightarrow \infty} \sup \{ |x| : x \in \text{spt}(\phi_m) \} = 0.$$

Show if $f \in L^p(\mathbb{R}^n)$, $\lim_{m \rightarrow \infty} f * \phi_m = f$ in $L^p(\mathbb{R}^n)$.

9. Let $\phi : \mathbb{R} \rightarrow \mathbb{R}$ be convex. This means $\phi(\lambda x + (1-\lambda)y) \leq \lambda \phi(x) + (1-\lambda)\phi(y)$ whenever $\lambda \in [0, 1]$. Verify that if $x < y < z$, then $\frac{\phi(y) - \phi(x)}{y-x} \leq \frac{\phi(z) - \phi(y)}{z-y}$ and that

$\frac{\phi(z)-\phi(x)}{z-x} \leq \frac{\phi(z)-\phi(y)}{z-y}$. Show if $s \in \mathbb{R}$ there exists λ such that $\phi(s) \leq \phi(t) + \lambda(s-t)$ for all t . Show that if ϕ is convex, then ϕ is continuous.

10. Let $\frac{1}{p} + \frac{1}{p'} = 1$, $p > 1$, let $f \in L^p(\mathbb{R})$, $g \in L^{p'}(\mathbb{R})$. Show $f * g$ is uniformly continuous on \mathbb{R} and $|(f * g)(x)| \leq \|f\|_{L^p} \|g\|_{L^{p'}}$. **Hint:** You need to consider why $f * g$ exists and then this follows from the definition of convolution and continuity of translation in L^p .
11. $B(p, q) = \int_0^1 x^{p-1} (1-x)^{q-1} dx$, $\Gamma(p) = \int_0^\infty e^{-t} t^{p-1} dt$ for $p, q > 0$. The first of these is called the beta function, while the second is the gamma function. Show a.) $\Gamma(p+1) = p\Gamma(p)$; b.) $\Gamma(p)\Gamma(q) = B(p, q)\Gamma(p+q)$.
12. Let $f \in C_c(0, \infty)$. Define $F(x) = \frac{1}{x} \int_0^x f(t) dt$. Show $\|F\|_{L^p(0, \infty)} \leq \frac{p}{p-1} \|f\|_{L^p(0, \infty)}$ whenever $p > 1$. **Hint:** Argue there is no loss of generality in assuming $f \geq 0$ and then assume this is so. Integrate $\int_0^\infty |F(x)|^p dx$ by parts as follows:

$$\int_0^\infty F^p dx = \overbrace{x F^p|_0^\infty}^{\text{show } = 0} - p \int_0^\infty x F^{p-1} F' dx.$$

Now show $x F' = f - F$ and use this in the last integral. Complete the argument by using Holder's inequality and $p-1 = p/q$.

13. \uparrow Now suppose $f \in L^p(0, \infty)$, $p > 1$, and f not necessarily in $C_c(0, \infty)$. Show that $F(x) = \frac{1}{x} \int_0^x f(t) dt$ still makes sense for each $x > 0$. Show the inequality of Problem 12 is still valid. This inequality is called Hardy's inequality. **Hint:** To show this, use the above inequality along with the density of $C_c(0, \infty)$ in $L^p(0, \infty)$.
14. Suppose $f, g \geq 0$. When does equality hold in Holder's inequality?
15. Show the Vitali Convergence theorem implies the Dominated Convergence theorem for finite measure spaces but there exist examples where the Vitali convergence theorem works and the dominated convergence theorem does not.
16. \uparrow Suppose $\mu(\Omega) < \infty$, $\{f_n\} \subseteq L^1(\Omega)$, and $\int_\Omega h(|f_n|) d\mu < C$ for all n where h is a continuous, nonnegative function satisfying $\lim_{t \rightarrow \infty} \frac{h(t)}{t} = \infty$. Show $\{f_n\}$ is uniformly integrable. In applications, this often occurs in the form of a bound on $\|f_n\|_p$.
17. \uparrow Sometimes, especially in books on probability, a different definition of uniform integrability is used than that presented here. A set of functions \mathfrak{S} , defined on a finite measure space, $(\Omega, \mathcal{S}, \mu)$ is said to be uniformly integrable if for all $\varepsilon > 0$ there exists $\alpha > 0$ such that for all $f \in \mathfrak{S}$, $\int_{|f| \geq \alpha} |f| d\mu \leq \varepsilon$. Show that this definition is equivalent to the definition of uniform integrability with the addition of the condition that there is a constant, $C < \infty$ such that $\int |f| d\mu \leq C$ for all $f \in \mathfrak{S}$.
18. Suppose $f \in L^\infty \cap L^1$. Show $\lim_{p \rightarrow \infty} \|f\|_{L^p} = \|f\|_\infty$. **Hint:**

$$(\|f\|_\infty - \varepsilon)^p \mu(|f| > \|f\|_\infty - \varepsilon) \leq \int_{|f| > \|f\|_\infty - \varepsilon} |f|^p d\mu \leq$$

$$\int |f|^p d\mu = \int |f|^{p-1} |f| d\mu \leq \|f\|_\infty^{p-1} \int |f| d\mu.$$

Now raise both ends to the $1/p$ power and take \liminf and \limsup as $p \rightarrow \infty$. You should get $\|f\|_\infty - \varepsilon \leq \liminf \|f\|_p \leq \limsup \|f\|_p \leq \|f\|_\infty$.

19. Suppose $\mu(\Omega) < \infty$. Show that if $1 \leq p < q$, then $L^q(\Omega) \subseteq L^p(\Omega)$. **Hint** Use Holder's inequality.
20. Show $L^1(\mathbb{R}) \not\subseteq L^2(\mathbb{R})$ and $L^2(\mathbb{R}) \not\subseteq L^1(\mathbb{R})$ if Lebesgue measure is used. **Hint:** Consider $1/\sqrt{x}$ and $1/x$.
21. Suppose that $\theta \in [0, 1]$ and $r, s, q > 0$ with $\frac{1}{q} = \frac{\theta}{r} + \frac{1-\theta}{s}$. Show that $(\int |f|^q d\mu)^{1/q} \leq (\int |f|^r d\mu)^{\theta/r} (\int |f|^s d\mu)^{(1-\theta)/s}$. If $q, r, s \geq 1$ this says that $\|f\|_q \leq \|f\|_r^\theta \|f\|_s^{1-\theta}$. Using this, show that $\ln(\|f\|_q) \leq \theta \ln(\|f\|_r) + (1-\theta) \ln(\|f\|_s)$. **Hint:** $\int |f|^q d\mu = \int |f|^{q\theta} |f|^{q(1-\theta)} d\mu$. Now note that $1 = \frac{\theta q}{r} + \frac{q(1-\theta)}{s}$ and use Holder's inequality.
22. Suppose f is a function in $L^1(\mathbb{R})$ and f is infinitely differentiable. Is $f' \in L^1(\mathbb{R})$? **Hint:** What if $\phi \in C_c^\infty(0, 1)$ and $f(x) = \phi(2^n(x-n))$ for $x \in (n, n+1)$, $f(x) = 0$ if $x < 0$?
23. Establish the following for $f \neq 0$ in L^p

$$\|f\|_p = \int_{\Omega} \frac{|f|^p}{\|f\|_p^{p-1}} d\mu = \int_{\Omega} f \frac{|f|^{p-2} \bar{f}}{\|f\|_p^{p-1}} \leq \sup_{\|g\|_q \leq 1} \int |f| |g| d\mu \leq \|f\|_p$$

24. ↑ From the above problem, if f is nonnegative and product measurable,

$$\left(\int \left(\int f(x, y) d\mu(x) \right)^p dv(y) \right)^{1/p} = \sup_{\|h\|_q \leq 1} \int \left(\int f(x, y) d\mu(x) \right) h(y) dv(y)$$

Now use Fubini's theorem and then the Holder inequality to obtain

$$= \sup_{\|h\|_q \leq 1} \int \int f(x, y) h(y) dv(y) d\mu(x) \leq \int \left(\int f(x, y)^p dv(y) \right)^{1/p} d\mu(x)$$

This gives another proof of the important Minkowski inequality for integrals.

25. Let $0 < p < 1$ and let f, g be measurable \mathbb{C} valued functions. Also

$$\int_{\Omega} |g|^{p/(p-1)} d\mu < \infty, \quad \int_{\Omega} |f|^p d\mu < \infty$$

Then show the following backwards Holder inequality holds.

$$\int_{\Omega} |fg| d\mu \geq \left(\int_{\Omega} |f|^p d\mu \right)^{1/p} \left(\int_{\Omega} |g|^{p/(p-1)} d\mu \right)^{(p-1)/p}$$

Hint: You should first note that $g = 0$ only on a set of measure zero. Then you could write $\int_{\Omega} |f|^p d\mu = \int_{\Omega} |g|^{-p} |fg|^p d\mu$. Apply the usual Holder inequality with $1/p$ one of the exponents and $1/(1-p)$ the other exponent. Then the above is $\leq \left(\int_{\Omega} |g|^{-p/(1-p)} d\mu \right)^{1-p} \left(\int_{\Omega} |fg|^p d\mu \right)^p$ etc. Note the usual Holder inequality in case $p > 1$ is $\int_{\Omega} |fg| d\mu \leq \left(\int_{\Omega} |f|^p d\mu \right)^{1/p} \left(\int_{\Omega} |g|^{p/(p-1)} d\mu \right)^{(p-1)/p}$.

26. Let $0 < p < 1$ and suppose $\int |h|^p d\mu < \infty$ for $h = f, g$. Then

$$\left(\int (|f| + |g|)^p d\mu \right)^{1/p} \geq \left(\int |f|^p d\mu \right)^{1/p} + \left(\int |g|^p d\mu \right)^{1/p}$$

This is the backwards Minkowski inequality. **Hint:** First explain why, since $p < 1$, $(|f| + |g|)^p \leq |f|^p + |g|^p$. It follows from this that

$$\int_{\Omega} \left((|f| + |g|)^{p-1} \right)^{p/(p-1)} d\mu < \infty$$

since $\int |h|^p d\mu < \infty$ for $h = f, g$. Then $(|f| + |g|)^{p-1}$ plays the role of $|g|$ in the above backwards Holder inequality. Next do this:

$$\begin{aligned} \int_{\Omega} (|f| + |g|)^p d\mu &= \int_{\Omega} (|f| + |g|)^{p-1} (|f| + |g|) d\mu \\ &= \int_{\Omega} (|f| + |g|)^{p-1} |f| d\mu + \int_{\Omega} (|f| + |g|)^{p-1} |g| d\mu \end{aligned}$$

Now apply the backwards Holder inequality of the above problem. The first term gives

$$\int_{\Omega} (|f| + |g|)^{p-1} |f| d\mu \geq \left(\int_{\Omega} (|f| + |g|)^p d\mu \right)^{(p-1)/p} \left(\int_{\Omega} |f|^p d\mu \right)^{1/p}$$

27. Let $f \in L^1_{loc}(\mathbb{R})$. Show there exists a set of measure zero N such that if $x \notin N$, then if $\{I_n\}$ is a sequence of intervals containing x such that $m(I_n) \rightarrow 0$ then

$$\frac{1}{m(I_n)} \int_{I_n} |f - f(x)| dx \rightarrow 0.$$

Generalize to higher dimensions if possible. Also, does I_n have to be an interval?

28. Suppose $F(x) = \int_a^x f(t) dt$ so that F is absolutely continuous where $f \in L^1([a, b])$. Show that $f \in L^p$ for $p > 1$ if and only if there is $M < \infty$ such that whenever $a = x_0 < x_1 < \dots < x_n = b$ it follows that $\sum_{i=1}^n \frac{|F(x_i) - F(x_{i-1})|^p}{(x_i - x_{i-1})^{p-1}} < M$. This problem is a result of F. Riesz. **Hint:** The first part is an easy application of Holder's inequality. For the second, let \mathcal{P}_n be a sequence of partitions of $[a, b]$ such that the subintervals have lengths converging to 0. Define $f_n(x) \equiv \sum_{k=1}^n \frac{F(x_k^n) - F(x_{k-1}^n)}{x_k^n - x_{k-1}^n} \chi_{I_k^n}(x)$ where the intervals of \mathcal{P}_n are $I_k^n = [x_{k-1}^n, x_k^n]$. Then for a.e. x , $f_n(x) \rightarrow f(x)$ thanks to the Lebesgue fundamental theorem of calculus and Problem 27. Now apply Fatou's lemma to say that $\int_a^b |f(x)|^p dx \leq \liminf_{n \rightarrow \infty} \int_a^b |f_n(x)|^p dx$ and simplify this last integral by breaking it into a sum of integrals over the sub-intervals of \mathcal{P}_n . Note $\frac{|F(x_k^n) - F(x_{k-1}^n)|^p}{(x_k^n - x_{k-1}^n)^p}$ does not depend on $x \in I_k^n$.
29. If $f \in L^p(U, m_p)$, where U is a bounded open set in \mathbb{R}^n , show there exists a sequence of smooth functions which converges to f a.e. Then show there exists a sequence of polynomials whose restriction to U converges a.e. to f on U .
30. In Corollary 12.4.3 can you generalize where f is only in $L^1_{loc}(\mathbb{R}^p, m_p)$.

Chapter 13

Fourier Transforms

13.1 Fourier Transforms of Functions in \mathcal{G}

First is a definition of a very specialized set of functions. Here the measure space will be $(\mathbb{R}^n, m_n, \mathcal{F}_n)$, familiar Lebesgue measure.

First recall the following definition of a polynomial.

Definition 13.1.1 $\alpha = (\alpha_1, \dots, \alpha_n)$ for $\alpha_1 \dots \alpha_n$ nonnegative integers is called a multi-index. For α a multi-index, $|\alpha| \equiv \alpha_1 + \dots + \alpha_n$ and if $\mathbf{x} \in \mathbb{R}^n, \mathbf{x} = (x_1, \dots, x_n)$, and f a function, define

$$\mathbf{x}^\alpha \equiv x_1^{\alpha_1} x_2^{\alpha_2} \dots x_n^{\alpha_n}.$$

A polynomial in n variables of degree m is a function of the form

$$p(\mathbf{x}) = \sum_{|\alpha| \leq m} a_\alpha \mathbf{x}^\alpha.$$

Here α is a multi-index as just described and $a_\alpha \in \mathbb{C}$. Also define for $\alpha = (\alpha_1, \dots, \alpha_n)$ a multi-index

$$D^\alpha f(\mathbf{x}) \equiv \frac{\partial^{|\alpha|} f}{\partial x_1^{\alpha_1} \partial x_2^{\alpha_2} \dots \partial x_n^{\alpha_n}}.$$

Definition 13.1.2 Define \mathcal{G}_1 to be the functions of the form $p(\mathbf{x}) e^{-a|\mathbf{x}|^2}$ where $a > 0$ is rational and $p(\mathbf{x})$ is a polynomial having all rational coefficients, a_α being “rational” if it is of the form $a + ib$ for $a, b \in \mathbb{Q}$. Let \mathcal{G} be all finite sums of functions in \mathcal{G}_1 . Thus \mathcal{G} is an algebra of functions which has the property that if $f \in \mathcal{G}$ then $\bar{f} \in \mathcal{G}$.

Thus there are countably many functions in \mathcal{G}_1 . This is because, for each m , there are countably many choices for a_α for $|\alpha| \leq m$ since there are finitely many α for $|\alpha| \leq m$ and for each such α , there are countably many choices for a_α since $\mathbb{Q} + i\mathbb{Q}$ is countable. (Why?) Thus there are countably many polynomials having degree no more than m . This is true for each m and so the number of different polynomials is a countable union of countable sets which is countable. Now there are countably many choices of $e^{-a|\mathbf{x}|^2}$ and so there are countably many in \mathcal{G}_1 because the Cartesian product of countable sets is countable.

Now \mathcal{G} consists of finite sums of functions in \mathcal{G}_1 . Therefore, it is countable because for each $m \in \mathbb{N}$, there are countably many such sums which are possible.

I will show now that \mathcal{G} is dense in $L^p(\mathbb{R}^n)$ but first, here is a lemma which follows from the Stone Weierstrass theorem.

Lemma 13.1.3 \mathcal{G} is dense in $C_0(\mathbb{R}^n)$ with respect to the norm,

$$\|f\|_\infty \equiv \sup \{|f(\mathbf{x})| : \mathbf{x} \in \mathbb{R}^n\}$$

Proof: By the Weierstrass approximation theorem, it suffices to show \mathcal{G} separates the points and annihilates no point. It was already observed in the above definition that $\bar{f} \in \mathcal{G}$ whenever $f \in \mathcal{G}$. If $\mathbf{y}_1 \neq \mathbf{y}_2$ suppose first that $|\mathbf{y}_1| \neq |\mathbf{y}_2|$. Then in this case, you can let $f(\mathbf{x}) \equiv e^{-|\mathbf{x}|^2}$. Then $f \in \mathcal{G}$ and $f(\mathbf{y}_1) \neq f(\mathbf{y}_2)$. If $|\mathbf{y}_1| = |\mathbf{y}_2|$, then suppose $y_{1k} \neq y_{2k}$. This must happen for some k because $\mathbf{y}_1 \neq \mathbf{y}_2$. Then let $f(\mathbf{x}) \equiv x_k e^{-|\mathbf{x}|^2}$. Thus \mathcal{G} separates points. Now $e^{-|\mathbf{x}|^2}$ is never equal to zero and so \mathcal{G} annihilates no point of \mathbb{R}^n . ■

These functions are clearly quite specialized. Therefore, the following theorem is somewhat surprising.

Theorem 13.1.4 *For each $p \geq 1, p < \infty, \mathcal{G}$ is dense in $L^p(\mathbb{R}^n)$. Since \mathcal{G} is countable, this shows that $L^p(\mathbb{R}^n)$ is separable.*

Proof: Let $f \in L^p(\mathbb{R}^n)$. Then there exists $g \in C_c(\mathbb{R}^n)$ such that $\|f - g\|_p < \varepsilon$. Now let $b > 0$ be large enough that

$$\int_{\mathbb{R}^n} \left(e^{-b|x|^2} \right)^p dx < \varepsilon^p.$$

Then $x \rightarrow g(x) e^{b|x|^2}$ is in $C_c(\mathbb{R}^n) \subseteq C_0(\mathbb{R}^n)$. Therefore, from Lemma 13.1.3 there exists $\psi \in \mathcal{G}$ such that

$$\left\| g e^{b|\cdot|^2} - \psi \right\|_\infty < 1$$

Therefore, letting $\phi(x) \equiv e^{-b|x|^2} \psi(x)$ it follows that $\phi \in \mathcal{G}$ and for all $x \in \mathbb{R}^n$,

$$|g(x) - \phi(x)| < e^{-b|x|^2}$$

Therefore, $(\int_{\mathbb{R}^n} |g(x) - \phi(x)|^p dx)^{1/p} \leq \left(\int_{\mathbb{R}^n} \left(e^{-b|x|^2} \right)^p dx \right)^{1/p} < \varepsilon$. It follows

$$\|f - \phi\|_p \leq \|f - g\|_p + \|g - \phi\|_p < 2\varepsilon. \blacksquare$$

From now on, we can drop the restriction that the coefficients of the polynomials in \mathcal{G} are rational. We also drop the restriction that a is rational. Thus \mathcal{G} will be finite sums of functions which are of the form $p(x) e^{-a|x|^2}$ where the coefficients of p are complex and $a > 0$.

The following lemma is also interesting even if it is obvious.

Lemma 13.1.5 *For $\psi \in \mathcal{G}$, p a polynomial, and α, β multi-indices, $D^\alpha \psi \in \mathcal{G}$ and $p\psi \in \mathcal{G}$. Also*

$$\sup\{|x^\beta D^\alpha \psi(x)| : x \in \mathbb{R}^n\} < \infty$$

Thus these special functions are infinitely differentiable (smooth). They also have the property that they and all their partial derivatives vanish as $|x| \rightarrow \infty$.

Let \mathcal{G} be the functions of Definition 13.1.2 except, for the sake of convenience, remove all references to rational numbers. Thus \mathcal{G} consists of finite sums of polynomials having coefficients in \mathbb{C} times $e^{-a|x|^2}$ for some $a > 0$. The idea is to first understand the Fourier transform on these very specialized functions.

Definition 13.1.6 *For $\psi \in \mathcal{G}$ Define the Fourier transform F and the inverse Fourier transform, F^{-1} by*

$$F\psi(t) \equiv (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{-it \cdot x} \psi(x) dx, \quad F^{-1}\psi(t) \equiv (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{it \cdot x} \psi(x) dx.$$

where $t \cdot x \equiv \sum_{i=1}^n t_i x_i$. Note there is no problem with this definition because ψ is in $L^1(\mathbb{R}^n)$ and therefore, $|e^{it \cdot x} \psi(x)| \leq |\psi(x)|$, an integrable function.

The following lemma follows from the dominated convergence theorem and routine manipulations.

Lemma 13.1.7 For $\psi \in \mathcal{G}$,

$$D_t^\alpha F(\psi) = (-i)^{|\alpha|} F(x^\alpha \psi(x)), \quad D_t^\alpha F^{-1}(\psi) = (i)^{|\alpha|} F^{-1}(x^\alpha \psi(x))$$

In this lemma, $x^\alpha \psi(x)$ denotes the function $x \rightarrow x^\alpha \psi(x)$.

One reason for using the functions \mathcal{G} is that it is very easy to compute the Fourier transform of these functions. The first thing to do is to verify F and F^{-1} map \mathcal{G} to \mathcal{G} and that $F^{-1} \circ F(\psi) = \psi$. Next is a simple lemma from calculus which is about the Fourier transforms of $e^{-c|t|^2}$.

Lemma 13.1.8 The following hold. ($c > 0$)

$$\begin{aligned} \left(\frac{1}{2\pi}\right)^{n/2} \int_{\mathbb{R}^n} e^{-c|t|^2} e^{-is \cdot t} dt &= \left(\frac{1}{2\pi}\right)^{n/2} \int_{\mathbb{R}^n} e^{-c|t|^2} e^{is \cdot t} dt \\ &= \left(\frac{1}{2\pi}\right)^{n/2} e^{-\frac{|s|^2}{4c}} \left(\frac{\sqrt{\pi}}{\sqrt{c}}\right)^n = \left(\frac{1}{2c}\right)^{n/2} e^{-\frac{1}{4c}|s|^2}. \end{aligned} \quad (13.1)$$

That is, $F(e^{-c|t|^2}) = \left(\frac{1}{2c}\right)^{n/2} e^{-\frac{1}{4c}|s|^2}$ and $F^{-1}\left(e^{-c|t|^2}\right) = \left(\frac{1}{2c}\right)^{n/2} e^{-\frac{1}{4c}|s|^2}$.

Proof: Consider first the case of one dimension. Let $H(s)$ be given by

$$H(s) \equiv \int_{\mathbb{R}} e^{-ct^2} e^{-ist} dt = \int_{\mathbb{R}} e^{-ct^2} \cos(st) dt$$

Then $H(0)^2 = \int_{\mathbb{R}} \int_{\mathbb{R}} e^{-c(t^2+s^2)} dt ds = \int_0^\infty \int_0^{2\pi} e^{-cr^2} r d\theta dr = \frac{1}{c} \pi$ so $H(0) = \sqrt{\frac{\pi}{c}}$. Then using the dominated convergence theorem to differentiate, $H'(s) + \frac{s}{2c} H(s) = 0$, $H(0) = \sqrt{\frac{\pi}{c}}$. Thus

$$\frac{d}{ds} \left(e^{s^2/4c} H(s) \right) = 0 \text{ and so } e^{s^2/4c} H(s) - \sqrt{\frac{\pi}{c}} = 0, \text{ so } H(s) = \sqrt{\frac{\pi}{c}} e^{-s^2/4c}$$

Hence

$$\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-ct^2} e^{-ist} dt = \sqrt{\frac{\pi}{c}} \frac{1}{\sqrt{2\pi}} e^{-\frac{s^2}{4c}} = \left(\frac{1}{2c}\right)^{1/2} e^{-\frac{s^2}{4c}}.$$

This proves the formula in the case of one dimension. The case of the inverse Fourier transform is similar. The n dimensional formula follows from Fubini's theorem. ■

With these formulas, it is easy to verify F, F^{-1} map \mathcal{G} to \mathcal{G} and $F \circ F^{-1} = F^{-1} \circ F = id$.

Theorem 13.1.9 Each of F and F^{-1} map \mathcal{G} to \mathcal{G} . Also $F^{-1} \circ F(\psi) = \psi$ and $F \circ F^{-1}(\psi) = \psi$.

Proof: It is obvious that F, F^{-1} map \mathcal{G} to \mathcal{G} using Lemmas 13.1.7, 13.1.8. Indeed, for $\psi(x) = x^\alpha e^{-c|x|^2}$,

$$\begin{aligned} F(\psi) &= (-i)^{|\alpha|} D^\alpha F\left(e^{-c|x|^2}\right) = (-i)^{|\alpha|} D^\alpha \left(\left(\frac{1}{2c}\right)^{n/2} e^{-\frac{1}{4c}|s|^2} \right) \\ &= (-i)^{|\alpha|} (-1)^{|\alpha|} \left(\frac{1}{2c}\right)^{|\alpha|+n/2} s^\alpha e^{-\frac{1}{4c}|s|^2} \end{aligned}$$

Similarly, for $\psi(x) = x^\alpha e^{-c|x|^2}$,

$$F^{-1}(\psi) = i^{|\alpha|} (-1)^{|\alpha|} \left(\frac{1}{2c}\right)^{|\alpha|+n/2} s^\alpha e^{-\frac{1}{4c}|s|^2} \quad (13.2)$$

For $\psi(x) = x^\alpha e^{-c|x|^2}$,

$$\begin{aligned} F^{-1} \circ F(\psi) &= F^{-1} \left((-i)^{|\alpha|} (-1)^{|\alpha|} \left(\frac{1}{2c}\right)^{|\alpha|+n/2} s^\alpha e^{-\frac{1}{4c}|s|^2} \right) \\ &= i^{|\alpha|} \left(\frac{1}{2c}\right)^{|\alpha|+n/2} F^{-1} \left(s^\alpha e^{-\frac{1}{4c}|s|^2} \right) \end{aligned}$$

From 13.2 with $c \rightarrow 1/(4c)$,

$$= i^{|\alpha|} \left(\frac{1}{2c}\right)^{|\alpha|+n/2} i^{|\alpha|} (-1)^{|\alpha|} \left(\frac{1}{2(1/4c)}\right)^{|\alpha|+n/2} s^\alpha e^{-\frac{1}{4(1/4c)}|s|^2} = s^\alpha e^{-c|s|^2}$$

Since \mathcal{G} consists of sums of multiples of such ψ , this has shown that $F^{-1} \circ F(\psi) = \psi$. ■

13.2 Fourier Transforms of Just about Anything

I will define Fourier Transforms of the linear functions acting on \mathcal{G} and then show that this includes virtually anything which could possibly be of any interest.

Definition 13.2.1 Let \mathcal{G}^* denote the vector space of linear functions defined on \mathcal{G} which have values in \mathbb{C} . Thus $T \in \mathcal{G}^*$ means $T : \mathcal{G} \rightarrow \mathbb{C}$ and T is linear,

$$T(a\psi + b\phi) = aT(\psi) + bT(\phi) \text{ for all } a, b \in \mathbb{C}, \psi, \phi \in \mathcal{G}$$

Let $\psi \in \mathcal{G}$. Then ψ is an element of \mathcal{G}^* by defining $\psi(\phi) \equiv \int_{\mathbb{R}^n} \psi(x) \phi(x) dx$.

Then we have the following important lemma.

Lemma 13.2.2 The following is obtained for all $\phi, \psi \in \mathcal{G}$.

$$F\psi(\phi) = \psi(F\phi), \quad F^{-1}\psi(\phi) = \psi(F^{-1}\phi)$$

Also if $\psi \in \mathcal{G}$ and $\psi = 0$ in \mathcal{G}^* so that $\psi(\phi) = 0$ for all $\phi \in \mathcal{G}$, then $\psi = 0$ as a function.

Proof:

$$\begin{aligned} F\psi(\phi) &\equiv \int_{\mathbb{R}^n} F\psi(t) \phi(t) dt = \int_{\mathbb{R}^n} \left(\frac{1}{2\pi}\right)^{n/2} \int_{\mathbb{R}^n} e^{-it \cdot x} \psi(x) dx \phi(t) dt \\ &= \int_{\mathbb{R}^n} \psi(x) \left(\frac{1}{2\pi}\right)^{n/2} \int_{\mathbb{R}^n} e^{-it \cdot x} \phi(t) dt dx = \int_{\mathbb{R}^n} \psi(x) F\phi(x) dx \equiv \psi(F\phi) \end{aligned}$$

The other claim is similar.

Suppose now $\psi(\phi) = 0$ for all $\phi \in \mathcal{G}$. Then $\int_{\mathbb{R}^n} \psi \phi dx = 0$ for all $\phi \in \mathcal{G}$. Therefore, this is true for $\phi = \bar{\psi}$ and so $\psi = 0$. ■

This lemma suggests a way to define the Fourier transform of something in \mathcal{G}^* .

Definition 13.2.3 For $T \in \mathcal{G}^*$, define $FT, F^{-1}T \in \mathcal{G}^*$ by

$$FT(\phi) \equiv T(F\phi), F^{-1}T(\phi) \equiv T(F^{-1}\phi)$$

Lemma 13.2.4 F and F^{-1} are both one to one, onto, and are inverses of each other.

Proof: First note F and F^{-1} are both linear. This follows directly from the definition. Suppose now $FT = 0$. Then $FT(\phi) = T(F\phi) = 0$ for all $\phi \in \mathcal{G}$. But F and F^{-1} map \mathcal{G} onto \mathcal{G} because if $\psi \in \mathcal{G}$, then as shown above, $\psi = F(F^{-1}(\psi))$. Therefore, $T = 0$ and so F is one to one. Similarly F^{-1} is one to one. Now

$$F^{-1}(FT)(\phi) \equiv (FT)(F^{-1}\phi) \equiv T(F(F^{-1}(\phi))) = T\phi.$$

Therefore, $F^{-1} \circ F(T) = T$. Similarly, $F \circ F^{-1}(T) = T$. Thus both F and F^{-1} are one to one and onto and are inverses of each other as suggested by the notation. ■

Probably the most interesting things in \mathcal{G}^* are functions of various kinds. The following lemma will be useful in considering this situation.

Lemma 13.2.5 If $f \in L^1_{loc}(\mathbb{R}^n)$ and $\int_{\mathbb{R}^n} f\phi dx = 0$ for all $\phi \in C_c(\mathbb{R}^n)$, then $f = 0$ a.e.

Proof: Let E be bounded and Lebesgue measurable. By regularity, there exists a compact set $K_k \subseteq E$ and an open set $V_k \supseteq E$ such that $m_n(V_k \setminus K_k) < 2^{-k}$. Let h_k equal 1 on K_k , vanish on V_k^C , and take values between 0 and 1. Let $E_k \equiv \left[|h_k - \mathcal{X}_{E_k}| > \left(\frac{2}{3}\right)^k\right]$. Then

$$m_n(E_k) \leq \left(\frac{3}{2}\right)^k \int_{E_k} |h_k - \mathcal{X}_{E_k}| dm_n < \left(\frac{3}{2}\right)^k \int_{V_k \setminus K_k} 2 dm_n = 2 \left(\frac{3}{2}\right)^k \frac{1}{2^k} = 2 \left(\frac{3}{4}\right)^k$$

and so $\sum_k m_n(E_k) < \infty$ so there is a set of measure zero such that off this set, $h_k \rightarrow \mathcal{X}_E$.

Hence, by the dominated convergence theorem, $\int f \mathcal{X}_E dm_n = \lim_{k \rightarrow \infty} \int f h_k dm_n = 0$. It follows that for E an arbitrary Lebesgue measurable set, $\int f \mathcal{X}_{B(0,R)} \mathcal{X}_E dm_n = 0$. Let

$$\operatorname{sgn} f = \begin{cases} \frac{\bar{f}}{|f|} & \text{if } |f| \neq 0 \\ 0 & \text{if } |f| = 0 \end{cases}$$

By Theorem 9.1.6 applied to positive and negative parts, there exists $\{s_k\}$, a sequence of simple functions converging pointwise to $\operatorname{sgn} f$ such that $|s_k| \leq 1$. Then by the dominated convergence theorem again, $\int |f| \mathcal{X}_{B(0,R)} dm_n = \lim_{k \rightarrow \infty} \int f \mathcal{X}_{B(0,R)} s_k dm_n = 0$. Since R is arbitrary, $|f| = 0$ a.e. ■

Corollary 13.2.6 Let $f \in L^1(\mathbb{R}^n)$ and suppose $\int_{\mathbb{R}^n} f(x)\phi(x) dx = 0$ for all $\phi \in \mathcal{G}$. Then $f = 0$ a.e.

Proof: Let $\psi \in C_c(\mathbb{R}^n)$. Then by the Stone Weierstrass approximation theorem, there exists a sequence of functions, $\{\phi_k\} \subseteq \mathcal{G}$ such that $\phi_k \rightarrow \psi$ uniformly. Then by the dominated convergence theorem, $\int f \psi dx = \lim_{k \rightarrow \infty} \int f \phi_k dx = 0$. By Lemma 13.2.5 $f = 0$. ■

The next theorem is the main result of this sort.

Theorem 13.2.7 Let $f \in L^p(\mathbb{R}^n)$, $p \geq 1$, or suppose f is measurable and has polynomial growth, $|f(x)| \leq K(1 + |x|^2)^m$ for some $m \in \mathbb{N}$. Then if $\int f \psi dx = 0$, for all $\psi \in \mathcal{G}$, then it follows $f = 0$.

Proof: First note that if $f \in L^p(\mathbb{R}^n)$ or has polynomial growth, then it makes sense to write the integral $\int f\psi dx$ described above. This is obvious in the case of polynomial growth. In the case where $f \in L^p(\mathbb{R}^n)$ it also makes sense because

$$\int |f||\psi| dx \leq \left(\int |f|^p dx \right)^{1/p} \left(\int |\psi|^{p'} dx \right)^{1/p'} < \infty$$

due to the fact mentioned above that all these functions in \mathcal{G} are in $L^p(\mathbb{R}^n)$ for every $p \geq 1$. Suppose now that $f \in L^p$, $p \geq 1$. The case where $f \in L^1(\mathbb{R}^n)$ was dealt with in Corollary 13.2.6. Suppose $f \in L^p(\mathbb{R}^n)$ for $p > 1$. Then

$$|f|^{p-2}\bar{f} \in L^{p'}(\mathbb{R}^n), \quad \left(p' = q, \frac{1}{p} + \frac{1}{q} = 1 \right)$$

and by density of \mathcal{G} in $L^{p'}(\mathbb{R}^n)$ (Theorem 13.1.4), there exists a sequence $\{g_k\} \subseteq \mathcal{G}$ such that $\|g_k - |f|^{p-2}\bar{f}\|_{p'} \rightarrow 0$. Then

$$\begin{aligned} \int_{\mathbb{R}^n} |f|^p dx &= \int_{\mathbb{R}^n} f(|f|^{p-2}\bar{f} - g_k) dx + \int_{\mathbb{R}^n} f g_k dx \\ &= \int_{\mathbb{R}^n} f(|f|^{p-2}\bar{f} - g_k) dx \leq \|f\|_{L^p} \|g_k - |f|^{p-2}\bar{f}\|_{p'} \end{aligned}$$

which converges to 0. Hence $f = 0$.

It remains to consider the case where f has polynomial growth. Thus $x \rightarrow f(x)e^{-|x|^2} \in L^1(\mathbb{R}^n)$. Therefore, for all $\psi \in \mathcal{G}$, $0 = \int f(x)e^{-|x|^2}\psi(x) dx$ because $e^{-|x|^2}\psi(x) \in \mathcal{G}$. Therefore, by the first part, $f(x)e^{-|x|^2} = 0$ a.e. ■

Note that “polynomial growth” could be replaced with a condition of the form

$$|f(x)| \leq K(1 + |x|^2)^m e^{k|x|^\alpha}, \quad \alpha < 2$$

and the same proof would yield that these functions are in \mathcal{G}^* . The main thing to observe is that almost all functions of interest are in \mathcal{G}^* .

Theorem 13.2.8 *Let f be a measurable function with polynomial growth,*

$$|f(x)| \leq C(1 + |x|^2)^N \text{ for some } N,$$

or let $f \in L^p(\mathbb{R}^n)$ for some $p \in [1, \infty]$. Then $f \in \mathcal{G}^$ if $f(\phi) \equiv \int f\phi dx$.*

Proof: Let f have polynomial growth first. Then the above integral is clearly well defined and so in this case, $f \in \mathcal{G}^*$.

Next suppose $f \in L^p(\mathbb{R}^n)$ with $\infty > p \geq 1$. Then it is clear again that the above integral is well defined because of the fact that ϕ is a sum of polynomials times exponentials of the form $e^{-c|x|^2}$ and these are in $L^{p'}(\mathbb{R}^n)$. Also $\phi \rightarrow f(\phi)$ is clearly linear in both cases. ■

This has shown that for nearly any reasonable function, you can define its Fourier transform as described above. You could also define the Fourier transform of a finite Borel measure μ because for such a measure $\psi \rightarrow \int_{\mathbb{R}^n} \psi d\mu$ is a linear functional on \mathcal{G} . This includes the very important case of probability distribution measures. The theoretical basis for this assertion will be given a little later.

13.2.1 Fourier Transforms of Functions in $L^1(\mathbb{R}^n)$

First suppose $f \in L^1(\mathbb{R}^n)$.

Theorem 13.2.9 *Let $f \in L^1(\mathbb{R}^n)$. Then $Ff(\phi) = \int_{\mathbb{R}^n} g\phi dt$ where*

$$g(t) = \left(\frac{1}{2\pi}\right)^{n/2} \int_{\mathbb{R}^n} e^{-it \cdot x} f(x) dx$$

and $F^{-1}f(\phi) = \int_{\mathbb{R}^n} g\phi dt$ where $g(t) = \left(\frac{1}{2\pi}\right)^{n/2} \int_{\mathbb{R}^n} e^{it \cdot x} f(x) dx$. In short,

$$Ff(t) \equiv (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{-it \cdot x} f(x) dx, \quad F^{-1}f(t) \equiv (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{it \cdot x} f(x) dx.$$

Proof: From the definition and Fubini's theorem,

$$\begin{aligned} Ff(\phi) &\equiv \int_{\mathbb{R}^n} f(t) F\phi(t) dt = \int_{\mathbb{R}^n} f(t) \left(\frac{1}{2\pi}\right)^{n/2} \int_{\mathbb{R}^n} e^{-it \cdot x} \phi(x) dx dt \\ &= \int_{\mathbb{R}^n} \left(\left(\frac{1}{2\pi}\right)^{n/2} \int_{\mathbb{R}^n} f(t) e^{-it \cdot x} dt \right) \phi(x) dx. \end{aligned}$$

Since $\phi \in \mathcal{G}$ is arbitrary, it follows from Theorem 13.2.7 that $Ff(x)$ is given by the claimed formula. The case of F^{-1} is identical. ■

Here are interesting properties of these Fourier transforms of functions in L^1 .

Theorem 13.2.10 *If $f \in L^1(\mathbb{R}^n)$ and $\|f_k - f\|_1 \rightarrow 0$, then Ff_k and $F^{-1}f_k$ converge uniformly to Ff and $F^{-1}f$ respectively. If $f \in L^1(\mathbb{R}^n)$, then $F^{-1}f$ and Ff are both continuous and bounded. Also,*

$$\lim_{|x| \rightarrow \infty} F^{-1}f(x) = \lim_{|x| \rightarrow \infty} Ff(x) = 0. \quad (13.3)$$

Furthermore, for $f \in L^1(\mathbb{R}^n)$ both Ff and $F^{-1}f$ are uniformly continuous.

Proof: The first claim follows from the following inequality.

$$\begin{aligned} |Ff_k(t) - Ff(t)| &\leq (2\pi)^{-n/2} \int_{\mathbb{R}^n} |e^{-it \cdot x} f_k(x) - e^{-it \cdot x} f(x)| dx \\ &= (2\pi)^{-n/2} \int_{\mathbb{R}^n} |f_k(x) - f(x)| dx = (2\pi)^{-n/2} \|f - f_k\|_1. \end{aligned}$$

which a similar argument holding for F^{-1} .

Now consider the second claim of the theorem.

$$|Ff(t) - Ff(t')| \leq (2\pi)^{-n/2} \int_{\mathbb{R}^n} |e^{-it \cdot x} - e^{-it' \cdot x}| |f(x)| dx$$

The integrand is bounded by $2|f(x)|$, a function in $L^1(\mathbb{R}^n)$ and converges to 0 as $t' \rightarrow t$ and so the dominated convergence theorem implies Ff is continuous. To see $Ff(t)$ is uniformly bounded,

$$|Ff(t)| \leq (2\pi)^{-n/2} \int_{\mathbb{R}^n} |f(x)| dx < \infty.$$

A similar argument gives the same conclusions for F^{-1} .

It remains to verify 13.3 and the claim that Ff and $F^{-1}f$ are uniformly continuous.

$$|Ff(t)| \leq \left| (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{-it \cdot x} f(x) dx \right|$$

Now let $\varepsilon > 0$ be given and let $g \in C_c^\infty(\mathbb{R}^n)$ such that $(2\pi)^{-n/2} \|g - f\|_1 < \varepsilon/2$. Then

$$\begin{aligned} |Ff(t)| &\leq (2\pi)^{-n/2} \int_{\mathbb{R}^n} |f(x) - g(x)| dx + \left| (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{-it \cdot x} g(x) dx \right| \\ &\leq \varepsilon/2 + \left| (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{-it \cdot x} g(x) dx \right|. \end{aligned}$$

Now integrating by parts, it follows that for $\|t\|_\infty \equiv \max\{|t_j| : j = 1, \dots, n\} > 0$

$$|Ff(t)| \leq \varepsilon/2 + (2\pi)^{-n/2} \left| \frac{1}{\|t\|_\infty} \int_{\mathbb{R}^n} \sum_{j=1}^n \left| \frac{\partial g(x)}{\partial x_j} \right| dx \right| \quad (13.4)$$

and this last expression converges to zero as $\|t\|_\infty \rightarrow \infty$. The reason for this is that if $t_j \neq 0$, integration by parts with respect to x_j gives

$$(2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{-it \cdot x} g(x) dx = (2\pi)^{-n/2} \frac{1}{-it_j} \int_{\mathbb{R}^n} e^{-it \cdot x} \frac{\partial g(x)}{\partial x_j} dx.$$

Therefore, choose the j for which $\|t\|_\infty = |t_j|$ and the result of 13.4 holds. Therefore, from 13.4, if $\|t\|_\infty$ is large enough, $|Ff(t)| < \varepsilon$. Similarly, $\lim_{\|t\|_\infty \rightarrow \infty} F^{-1}(t) = 0$.

Consider the claim about uniform continuity. Let $\varepsilon > 0$ be given. Then there exists R such that if $\|t\|_\infty > R$, then $|Ff(t)| < \frac{\varepsilon}{2}$. Since Ff is continuous, it is uniformly continuous on the compact set $[-R-1, R+1]^n$. Therefore, there exists δ_1 such that if $\|t - t'\|_\infty < \delta_1$ for $t', t \in [-R-1, R+1]^n$, then

$$|Ff(t) - Ff(t')| < \varepsilon/2. \quad (13.5)$$

Now let $0 < \delta < \min(\delta_1, 1)$ and suppose $\|t - t'\|_\infty < \delta$. If both t, t' are contained in $[-R, R]^n$, then 13.5 holds. If $t \in [-R, R]^n$ and $t' \notin [-R, R]^n$, then both are contained in $[-R-1, R+1]^n$ and so this verifies 13.5 in this case. The other case is that neither point is in $[-R, R]^n$ and in this case,

$$|Ff(t) - Ff(t')| \leq |Ff(t)| + |Ff(t')| < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon. \quad \blacksquare$$

There is a very interesting relation between the Fourier transform and convolutions.

Theorem 13.2.11 *Let $f, g \in L^1(\mathbb{R}^n)$. Then $f * g \in L^1$, $F(f * g) = (2\pi)^{n/2} Ff Fg$.*

Proof: Consider

$$\int_{\mathbb{R}^n} \int_{\mathbb{R}^n} |f(x - y) g(y)| dy dx.$$

The function, $(x, y) \rightarrow |f(x - y) g(y)|$ is Lebesgue measurable and so by Fubini's theorem,

$$\int_{\mathbb{R}^n} \int_{\mathbb{R}^n} |f(x - y) g(y)| dy dx = \int_{\mathbb{R}^n} \int_{\mathbb{R}^n} |f(x - y) g(y)| dx dy = \|f\|_1 \|g\|_1 < \infty.$$

It follows that for a.e. \mathbf{x} , $\int_{\mathbb{R}^n} |f(\mathbf{x} - \mathbf{y}) g(\mathbf{y})| d\mathbf{y} < \infty$ and for each of these values of \mathbf{x} , it follows that $\int_{\mathbb{R}^n} f(\mathbf{x} - \mathbf{y}) g(\mathbf{y}) d\mathbf{y}$ exists and equals a function of \mathbf{x} which is in $L^1(\mathbb{R}^n)$, $f * g(\mathbf{x})$. Now

$$\begin{aligned} F(f * g)(\mathbf{t}) &\equiv (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{-it \cdot \mathbf{x}} f * g(\mathbf{x}) d\mathbf{x} \\ &= (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{-it \cdot \mathbf{x}} \int_{\mathbb{R}^n} f(\mathbf{x} - \mathbf{y}) g(\mathbf{y}) d\mathbf{y} d\mathbf{x} \\ &= (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{-it \cdot \mathbf{y}} g(\mathbf{y}) \int_{\mathbb{R}^n} e^{-it \cdot (\mathbf{x} - \mathbf{y})} f(\mathbf{x} - \mathbf{y}) d\mathbf{x} d\mathbf{y} \\ &= (2\pi)^{n/2} Ff(\mathbf{t}) Fg(\mathbf{t}). \blacksquare \end{aligned}$$

There are many other considerations involving Fourier transforms of functions in L^1 . Some others are in the exercises.

13.2.2 Fourier Transforms of Functions in $L^2(\mathbb{R}^n)$

Consider Ff and $F^{-1}f$ for $f \in L^2(\mathbb{R}^n)$. First note that the formula given for Ff and $F^{-1}f$ when $f \in L^1(\mathbb{R}^n)$ will not work for $f \in L^2(\mathbb{R}^n)$ unless f is also in $L^1(\mathbb{R}^n)$. Recall that $\overline{a + ib} = a - ib$.

Theorem 13.2.12 For $\phi \in \mathcal{G}$, $\|F\phi\|_2 = \|F^{-1}\phi\|_2 = \|\phi\|_2$.

Proof: First note that for $\psi \in \mathcal{G}$,

$$F(\overline{\psi}) = \overline{F^{-1}(\psi)}, \quad F^{-1}(\overline{\psi}) = \overline{F(\psi)}. \quad (13.6)$$

This follows from the definition. For example,

$$F\overline{\psi}(\mathbf{t}) = (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{-it \cdot \mathbf{x}} \overline{\psi}(\mathbf{x}) d\mathbf{x} = \overline{(2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{it \cdot \mathbf{x}} \psi(\mathbf{x}) d\mathbf{x}}$$

Let $\phi, \psi \in \mathcal{G}$. It was shown above that $\int_{\mathbb{R}^n} (F\phi)\psi(\mathbf{t}) d\mathbf{t} = \int_{\mathbb{R}^n} \phi(F\psi) d\mathbf{x}$. Similarly,

$$\int_{\mathbb{R}^n} \phi(F^{-1}\psi) d\mathbf{x} = \int_{\mathbb{R}^n} (F^{-1}\phi)\psi d\mathbf{t}. \quad (13.7)$$

Now, 13.6 - 13.7 imply

$$\begin{aligned} \int_{\mathbb{R}^n} |\phi|^2 d\mathbf{x} &= \int_{\mathbb{R}^n} \phi \overline{\phi} d\mathbf{x} = \int_{\mathbb{R}^n} \phi (\overline{F^{-1}(F\phi)}) d\mathbf{x} = \int_{\mathbb{R}^n} \phi F(\overline{F\phi}) d\mathbf{x} \\ &= \int_{\mathbb{R}^n} F\phi (\overline{F\phi}) d\mathbf{x} = \int_{\mathbb{R}^n} |F\phi|^2 d\mathbf{x}. \end{aligned}$$

Similarly $\|\phi\|_2 = \|F^{-1}\phi\|_2$. \blacksquare

Lemma 13.2.13 Let $f \in L^2(\mathbb{R}^n)$ and let $\phi_k \rightarrow f$ in $L^2(\mathbb{R}^n)$ where $\phi_k \in \mathcal{G}$. (Such a sequence exists because of density of \mathcal{G} in $L^2(\mathbb{R}^n)$.) Then Ff and $F^{-1}f$ are both in $L^2(\mathbb{R}^n)$ and the following limits take place in L^2 .

$$\lim_{k \rightarrow \infty} F(\phi_k) = F(f), \quad \lim_{k \rightarrow \infty} F^{-1}(\phi_k) = F^{-1}(f).$$

Proof: Let $\psi \in \mathcal{G}$ be given. Then from Theorem 13.2.8,

$$\begin{aligned} Ff(\psi) &\equiv f(F\psi) \equiv \int_{\mathbb{R}^n} f(x) F\psi(x) dx \\ &= \lim_{k \rightarrow \infty} \int_{\mathbb{R}^n} \phi_k(x) F\psi(x) dx = \lim_{k \rightarrow \infty} \int_{\mathbb{R}^n} F\phi_k(x) \psi(x) dx. \end{aligned}$$

Also by Theorem 13.2.12 $\{F\phi_k\}_{k=1}^\infty$ is Cauchy in $L^2(\mathbb{R}^n)$ since

$$\|F\phi_k - F\phi_l\|_{L^2} = \|\phi_k - \phi_l\|_{L^2},$$

and so $\lim_{k \rightarrow \infty} F\phi_k = h$ for some $h \in L^2(\mathbb{R}^n)$. Therefore, from the above, $Ff(\psi) = \int_{\mathbb{R}^n} h(x) \psi(x) dx$ which shows that $F(f) \in L^2(\mathbb{R}^n)$ and $h = F(f)$. The case of F^{-1} is entirely similar. ■

Since Ff and $F^{-1}f$ are in $L^2(\mathbb{R}^n)$, this also proves the following theorem.

Theorem 13.2.14 *If $f \in L^2(\mathbb{R}^n)$, Ff and $F^{-1}f$ are the unique elements of $L^2(\mathbb{R}^n)$ such that for all $\phi \in \mathcal{G}$,*

$$\int_{\mathbb{R}^n} Ff(x) \phi(x) dx = \int_{\mathbb{R}^n} f(x) F\phi(x) dx, \quad (13.8)$$

$$\int_{\mathbb{R}^n} F^{-1}f(x) \phi(x) dx = \int_{\mathbb{R}^n} f(x) F^{-1}\phi(x) dx. \quad (13.9)$$

Theorem 13.2.15 (Plancherel)

$$\|f\|_2 = \|Ff\|_2 = \|F^{-1}f\|_2. \quad (13.10)$$

Proof: Use the density of \mathcal{G} in $L^2(\mathbb{R}^n)$ to obtain a sequence, $\{\phi_k\}$ converging to f in $L^2(\mathbb{R}^n)$. Then by Lemma 13.2.13

$$\|Ff\|_2 = \lim_{k \rightarrow \infty} \|F\phi_k\|_2 = \lim_{k \rightarrow \infty} \|\phi_k\|_2 = \|f\|_2.$$

Similarly, $\|f\|_2 = \|F^{-1}f\|_2$. ■

The following corollary is a generalization of this. To prove this corollary, use the following simple lemma which comes as a consequence of the Cauchy Schwarz inequality.

Lemma 13.2.16 *Suppose $f_k \rightarrow f$ in $L^2(\mathbb{R}^n)$ and $g_k \rightarrow g$ in $L^2(\mathbb{R}^n)$. Then*

$$\lim_{k \rightarrow \infty} \int_{\mathbb{R}^n} f_k g_k dx = \int_{\mathbb{R}^n} f g dx$$

Proof:

$$\begin{aligned} \left| \int_{\mathbb{R}^n} f_k g_k dx - \int_{\mathbb{R}^n} f g dx \right| &\leq \left| \int_{\mathbb{R}^n} f_k g_k dx - \int_{\mathbb{R}^n} f_k g dx \right| + \left| \int_{\mathbb{R}^n} f_k g dx - \int_{\mathbb{R}^n} f g dx \right| \\ &\leq \|f_k\|_2 \|g - g_k\|_2 + \|g\|_2 \|f_k - f\|_2. \end{aligned}$$

Now $\|f_k\|_2$ is a Cauchy sequence and so it is bounded independent of k . Therefore, the above expression is smaller than ε whenever k is large enough. ■

Corollary 13.2.17 For $f, g \in L^2(\mathbb{R}^n)$,

$$\int_{\mathbb{R}^n} f \bar{g} dx = \int_{\mathbb{R}^n} Ff \overline{Fg} dx = \int_{\mathbb{R}^n} F^{-1}f \overline{F^{-1}g} dx.$$

Proof: First note the above formula is obvious if $f, g \in \mathcal{G}$. To see this, note

$$\begin{aligned} \int_{\mathbb{R}^n} Ff \overline{Fg} dx &= \int_{\mathbb{R}^n} Ff(x) \overline{\frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} e^{-i\mathbf{x} \cdot \mathbf{t}} g(\mathbf{t}) dt} dx \\ &= \int_{\mathbb{R}^n} \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} e^{i\mathbf{x} \cdot \mathbf{t}} Ff(x) dx \overline{g(\mathbf{t})} dt = \int_{\mathbb{R}^n} (F^{-1} \circ F) f(\mathbf{t}) \overline{g(\mathbf{t})} dt = \int_{\mathbb{R}^n} f(\mathbf{t}) \overline{g(\mathbf{t})} dt. \end{aligned}$$

The formula with F^{-1} is exactly similar.

Now to verify the corollary, let $\phi_k \rightarrow f$ in $L^2(\mathbb{R}^n)$ and let $\psi_k \rightarrow g$ in $L^2(\mathbb{R}^n)$. Then by Lemma 13.2.13

$$\int_{\mathbb{R}^n} Ff \overline{Fg} dx = \lim_{k \rightarrow \infty} \int_{\mathbb{R}^n} F\phi_k \overline{F\psi_k} dx = \lim_{k \rightarrow \infty} \int_{\mathbb{R}^n} \phi_k \overline{\psi_k} dx = \int_{\mathbb{R}^n} f \bar{g} dx$$

A similar argument holds for F^{-1} . ■

How does one compute Ff and $F^{-1}f$?

Theorem 13.2.18 For $f \in L^2(\mathbb{R}^n)$, let $f_r = f \chi_{E_r}$ where E_r is a bounded measurable set with $E_r \uparrow \mathbb{R}^n$. Then the following limits hold in $L^2(\mathbb{R}^n)$.

$$Ff = \lim_{r \rightarrow \infty} Ff_r, \quad F^{-1}f = \lim_{r \rightarrow \infty} F^{-1}f_r.$$

Proof: $\|f - f_r\|_2 \rightarrow 0$ and so $\|Ff - Ff_r\|_2 \rightarrow 0$ and $\|F^{-1}f - F^{-1}f_r\|_2 \rightarrow 0$ which both follow from Plancherel's Theorem. ■

What are Ff_r and $F^{-1}f_r$? Let $\phi \in \mathcal{G}$

$$\begin{aligned} \int_{\mathbb{R}^n} Ff_r \phi dx &= \int_{\mathbb{R}^n} f_r F\phi dx = (2\pi)^{-\frac{n}{2}} \int_{\mathbb{R}^n} \int_{\mathbb{R}^n} f_r(\mathbf{x}) e^{-i\mathbf{x} \cdot \mathbf{y}} \phi(\mathbf{y}) dy dx \\ &= \int_{\mathbb{R}^n} [(2\pi)^{-\frac{n}{2}} \int_{\mathbb{R}^n} f_r(\mathbf{x}) e^{-i\mathbf{x} \cdot \mathbf{y}} dx] \phi(\mathbf{y}) dy. \end{aligned}$$

Since this holds for all $\phi \in \mathcal{G}$, a dense subset of $L^2(\mathbb{R}^n)$, it follows that

$$Ff_r(\mathbf{y}) = (2\pi)^{-\frac{n}{2}} \int_{\mathbb{R}^n} f_r(\mathbf{x}) e^{-i\mathbf{x} \cdot \mathbf{y}} dx.$$

Similarly

$$F^{-1}f_r(\mathbf{y}) = (2\pi)^{-\frac{n}{2}} \int_{\mathbb{R}^n} f_r(\mathbf{x}) e^{i\mathbf{x} \cdot \mathbf{y}} dx.$$

This shows that to take the Fourier transform of a function in $L^2(\mathbb{R}^n)$, it suffices to take the limit as $r \rightarrow \infty$ in $L^2(\mathbb{R}^n)$ of $(2\pi)^{-\frac{n}{2}} \int_{\mathbb{R}^n} f_r(\mathbf{x}) e^{-i\mathbf{x} \cdot \mathbf{y}} dx$. A similar procedure works for the inverse Fourier transform.

Note this reduces to the earlier definition in case $f \in L^1(\mathbb{R}^n)$. Now consider the convolution of a function in L^2 with one in L^1 .

Theorem 13.2.19 *Let $h \in L^2(\mathbb{R}^n)$ and let $f \in L^1(\mathbb{R}^n)$. Then $h * f \in L^2(\mathbb{R}^n)$,*

$$F^{-1}(h * f) = (2\pi)^{n/2} F^{-1} h F^{-1} f, \quad F(h * f) = (2\pi)^{n/2} F h F f,$$

and

$$\|h * f\|_2 \leq \|h\|_2 \|f\|_1. \quad (13.11)$$

Proof: An application of Minkowski's inequality yields

$$\left(\int_{\mathbb{R}^n} \left(\int_{\mathbb{R}^n} |h(\mathbf{x} - \mathbf{y})| |f(\mathbf{y})| d\mathbf{y} \right)^2 d\mathbf{x} \right)^{1/2} \leq \|f\|_1 \|h\|_2. \quad (13.12)$$

Hence $\int |h(\mathbf{x} - \mathbf{y})| |f(\mathbf{y})| d\mathbf{y} < \infty$ a.e. \mathbf{x} and $\mathbf{x} \rightarrow \int h(\mathbf{x} - \mathbf{y}) f(\mathbf{y}) d\mathbf{y}$ is in $L^2(\mathbb{R}^n)$. Let $E_r \uparrow \mathbb{R}^n$, $m(E_r) < \infty$. Thus, $h_r \equiv \chi_{E_r} h \in L^2(\mathbb{R}^n) \cap L^1(\mathbb{R}^n)$, and letting $\phi \in \mathcal{G}$,

$$\begin{aligned} & \int F(h_r * f)(\phi) d\mathbf{x} = \\ & \equiv \int (h_r * f)(F\phi) d\mathbf{x} = (2\pi)^{-n/2} \int \int \int h_r(\mathbf{x} - \mathbf{y}) f(\mathbf{y}) e^{-i\mathbf{x} \cdot \mathbf{t}} \phi(\mathbf{t}) dt d\mathbf{y} d\mathbf{x} \\ & = (2\pi)^{-n/2} \int \int \left(\int h_r(\mathbf{x} - \mathbf{y}) e^{-i(\mathbf{x} - \mathbf{y}) \cdot \mathbf{t}} d\mathbf{x} \right) f(\mathbf{y}) e^{-i\mathbf{y} \cdot \mathbf{t}} d\mathbf{y} \phi(\mathbf{t}) dt \\ & = \int (2\pi)^{n/2} F h_r(\mathbf{t}) F f(\mathbf{t}) \phi(\mathbf{t}) dt. \end{aligned}$$

Since ϕ is arbitrary and \mathcal{G} is dense in $L^2(\mathbb{R}^n)$, $F(h_r * f) = (2\pi)^{n/2} F h_r F f$. Now by Minkowski's Inequality, $h_r * f \rightarrow h * f$ in $L^2(\mathbb{R}^n)$ and also it is clear that $h_r \rightarrow h$ in $L^2(\mathbb{R}^n)$; so, by Plancherel's theorem, you may take the limit in the above and conclude the following equation: $F(h * f) = (2\pi)^{n/2} F h F f$. The assertion for F^{-1} is similar and 13.11 follows from 13.12. ■

13.2.3 The Schwartz Class

The problem with \mathcal{G} is that it does not contain $C_c^\infty(\mathbb{R}^n)$. I have used it in presenting the Fourier transform because the functions in \mathcal{G} have a very specific form which made some technical details work out easier than in any other approach I have seen. The Schwartz class is a larger class of functions which does contain $C_c^\infty(\mathbb{R}^n)$ and also has the same nice properties as \mathcal{G} . The functions in the Schwartz class are infinitely differentiable and they vanish very rapidly as $|\mathbf{x}| \rightarrow \infty$ along with all their partial derivatives. This is the description of these functions, not a specific form involving polynomials times $e^{-\alpha|\mathbf{x}|^2}$. To describe this precisely requires some notation.

Definition 13.2.20 *$f \in \mathfrak{S}$, the Schwartz class, if $f \in C^\infty(\mathbb{R}^n)$ and for all positive integers N , $\rho_N(f) < \infty$ where*

$$\rho_N(f) = \sup\{(1 + |\mathbf{x}|^2)^N |D^\alpha f(\mathbf{x})| : \mathbf{x} \in \mathbb{R}^n, |\alpha| \leq N\}.$$

Thus $f \in \mathfrak{S}$ if and only if $f \in C^\infty(\mathbb{R}^n)$ and

$$\sup\{|\mathbf{x}^\beta D^\alpha f(\mathbf{x})| : \mathbf{x} \in \mathbb{R}^n\} < \infty \quad (13.13)$$

for all multi indices α and β .

Also note that if $f \in \mathfrak{S}$, then $p(f) \in \mathfrak{S}$ for any polynomial, p with $p(0) = 0$ and that

$$\mathfrak{S} \subseteq L^p(\mathbb{R}^n) \cap L^\infty(\mathbb{R}^n)$$

for any $p \geq 1$. To see this assertion about the $p(f)$, it suffices to consider the case of the product of two elements of the Schwartz class. If $f, g \in \mathfrak{S}$, then $D^\alpha(fg)$ is a finite sum of derivatives of f times derivatives of g . Therefore, $\rho_N(fg) < \infty$ for all N . You may wonder about examples of things in \mathfrak{S} . Clearly any function in $C_c^\infty(\mathbb{R}^n)$ is in \mathfrak{S} . However there are other functions in \mathfrak{S} . For example $e^{-|x|^2}$ is in \mathfrak{S} as you can verify for yourself and so is any function from \mathcal{G} . Note also that the density of $C_c(\mathbb{R}^n)$ in $L^p(\mathbb{R}^n)$ shows that \mathfrak{S} is dense in $L^p(\mathbb{R}^n)$ for every p .

Recall the Fourier transform of a function in $L^1(\mathbb{R}^n)$ is given by

$$Ff(t) \equiv (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{-it \cdot x} f(x) dx.$$

Therefore, this gives the Fourier transform for $f \in \mathfrak{S}$. The nice property which \mathfrak{S} has in common with \mathcal{G} is that the Fourier transform and its inverse map \mathfrak{S} one to one onto \mathfrak{S} . This means I could have presented the whole of the above theory in terms of \mathfrak{S} rather than in terms of \mathcal{G} . However, it is more technical.

Theorem 13.2.21 *If $f \in \mathfrak{S}$, then Ff and $F^{-1}f$ are also in \mathfrak{S} .*

Proof: To begin with, let $\alpha = e_j = (0, 0, \dots, 1, 0, \dots, 0)$, the 1 in the j^{th} slot.

$$\frac{F^{-1}f(t + he_j) - F^{-1}f(t)}{h} = (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{it \cdot x} f(x) \left(\frac{e^{ihx_j} - 1}{h} \right) dx. \quad (13.14)$$

Consider the integrand in 13.14.

$$\begin{aligned} \left| e^{it \cdot x} f(x) \left(\frac{e^{ihx_j} - 1}{h} \right) \right| &= |f(x)| \left| \left(\frac{e^{i(h/2)x_j} - e^{-i(h/2)x_j}}{h} \right) \right| \\ &= |f(x)| \left| \frac{i \sin((h/2)x_j)}{(h/2)} \right| \leq |f(x)| |x_j| \end{aligned}$$

and this is a function in $L^1(\mathbb{R}^n)$ because $f \in \mathfrak{S}$. Therefore by the Dominated Convergence Theorem,

$$\frac{\partial F^{-1}f(t)}{\partial t_j} = (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{it \cdot x} i x_j f(x) dx = i (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{it \cdot x} x^j f(x) dx.$$

Now $x^j f(x) \in \mathfrak{S}$ and so one can continue in this way and take derivatives indefinitely. Thus $F^{-1}f \in C^\infty(\mathbb{R}^n)$ and from the above argument,

$$D^\alpha F^{-1}f(t) = (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{it \cdot x} (ix)^\alpha f(x) dx.$$

To complete showing $F^{-1}f \in \mathfrak{S}$,

$$t^\beta D^\alpha F^{-1}f(t) = (2\pi)^{-n/2} \int_{\mathbb{R}^n} e^{it \cdot x} t^\beta (ix)^\alpha f(x) dx.$$

Integrate this integral by parts to get

$$t^\beta D^\alpha F^{-1} f(t) = (2\pi)^{-n/2} \int_{\mathbb{R}^n} i^{|\beta|} e^{it \cdot x} D^\beta ((ix)^\alpha f(x)) dx. \quad (13.15)$$

Here is how this is done.

$$\begin{aligned} \int_{\mathbb{R}} e^{it_j x_j} t_j^{\beta_j} (ix)^\alpha f(x) dx_j &= \frac{e^{it_j x_j} t_j^{\beta_j} (ix)^\alpha f(x)}{it_j} \Big|_{-\infty}^{\infty} + \\ &+ i \int_{\mathbb{R}} e^{it_j x_j} t_j^{\beta_j-1} D^{e_j} ((ix)^\alpha f(x)) dx_j \end{aligned}$$

where the boundary term vanishes because $f \in \mathfrak{S}$. Returning to 13.15, use the fact that $|e^{ia}| = 1$ to conclude

$$\left| t^\beta D^\alpha F^{-1} f(t) \right| \leq C \int_{\mathbb{R}^n} \left| D^\beta ((ix)^\alpha f(x)) \right| dx < \infty.$$

It follows $F^{-1} f \in \mathfrak{S}$. Similarly $F f \in \mathfrak{S}$ whenever $f \in \mathfrak{S}$. ■

Of course \mathfrak{S} can be considered a subset of \mathcal{G}^* as follows. For $\psi \in \mathfrak{S}$, $\psi(\phi) \equiv \int_{\mathbb{R}^n} \psi \phi dx$.

Theorem 13.2.22 *Let $\psi \in \mathfrak{S}$. Then $(F \circ F^{-1})(\psi) = \psi$ and $(F^{-1} \circ F)(\psi) = \psi$ whenever $\psi \in \mathfrak{S}$. Also F and F^{-1} map \mathfrak{S} one to one and onto \mathfrak{S} .*

Proof: The first claim follows from the fact that F and F^{-1} are inverses of each other on \mathcal{G}^* which was established above. For the second, let $\psi \in \mathfrak{S}$. Then $\psi = F(F^{-1}\psi)$. Thus F maps \mathfrak{S} onto \mathfrak{S} . If $F\psi = 0$, then do F^{-1} to both sides to conclude $\psi = 0$. Thus F is one to one and onto. Similarly, F^{-1} is one to one and onto. ■

13.2.4 Convolution

To begin with it is necessary to discuss the meaning of ϕf where $f \in \mathcal{G}^*$ and $\phi \in \mathcal{G}$. What should it mean? First suppose $f \in L^p(\mathbb{R}^n)$ or measurable with polynomial growth. Then ϕf also has these properties. Hence, it should be the case that $\phi f(\psi) = \int_{\mathbb{R}^n} \phi f \psi dx = \int_{\mathbb{R}^n} f(\phi \psi) dx$. This motivates the following definition.

Definition 13.2.23 *Let $T \in \mathcal{G}^*$ and let $\phi \in \mathcal{G}$. Then $\phi T \equiv T\phi \in \mathcal{G}^*$ will be defined by*

$$\phi T(\psi) \equiv T(\phi \psi).$$

The next topic is that of convolution. It was just shown that

$$F(f * \phi) = (2\pi)^{n/2} F\phi Ff, \quad F^{-1}(f * \phi) = (2\pi)^{n/2} F^{-1}\phi F^{-1}f$$

whenever $f \in L^2(\mathbb{R}^n)$ and $\phi \in \mathcal{G}$ so the same definition is retained in the general case because it makes perfect sense and agrees with the earlier definition.

Definition 13.2.24 *Let $f \in \mathcal{G}^*$ and let $\phi \in \mathcal{G}$. Then define the convolution of f with an element of \mathcal{G} as follows.*

$$f * \phi \equiv (2\pi)^{n/2} F^{-1}(F\phi Ff) \in \mathcal{G}^*$$

There is an obvious question. With this definition, is it true that

$$F^{-1}(f * \phi) = (2\pi)^{n/2} F^{-1} \phi F^{-1} f$$

as it was earlier?

Theorem 13.2.25 *Let $f \in \mathcal{G}^*$ and let $\phi \in \mathcal{G}$.*

$$F(f * \phi) = (2\pi)^{n/2} F\phi Ff, \quad (13.16)$$

$$F^{-1}(f * \phi) = (2\pi)^{n/2} F^{-1} \phi F^{-1} f. \quad (13.17)$$

Proof: Note that 13.16 follows from Definition 13.2.24 and both assertions hold for $f \in \mathcal{G}$. Consider 13.17. Here is a simple formula involving a pair of functions in \mathcal{G} .

$$\begin{aligned} & (\psi * F^{-1} F^{-1} \phi)(x) \\ &= \left(\int \int \int \psi(x - y) e^{iy \cdot y_1} e^{iy_1 \cdot z} \phi(z) dz dy_1 dy \right) (2\pi)^n \\ &= \left(\int \int \int \psi(x - y) e^{-iy \cdot \tilde{y}_1} e^{-i\tilde{y}_1 \cdot z} \phi(z) dz d\tilde{y}_1 dy \right) (2\pi)^n = (\psi * FF\phi)(x). \end{aligned}$$

Now for $\psi \in \mathcal{G}$,

$$\begin{aligned} (2\pi)^{n/2} F(F^{-1} \phi F^{-1} f)(\psi) &\equiv (2\pi)^{n/2} (F^{-1} \phi F^{-1} f)(F\psi) \equiv \\ (2\pi)^{n/2} F^{-1} f(F^{-1} \phi F\psi) &\equiv (2\pi)^{n/2} f(F^{-1} (F^{-1} \phi F\psi)) = \\ f\left((2\pi)^{n/2} F^{-1} ((FF^{-1} F^{-1} \phi)(F\psi))\right) &\equiv \\ f(\psi * F^{-1} F^{-1} \phi) &= f(\psi * FF\phi) \end{aligned} \quad (13.18)$$

Also

$$\begin{aligned} (2\pi)^{n/2} F^{-1}(F\phi Ff)(\psi) &\equiv (2\pi)^{n/2} (F\phi Ff)(F^{-1}\psi) \equiv \\ (2\pi)^{n/2} Ff(F\phi F^{-1}\psi) &\equiv (2\pi)^{n/2} f(F(F\phi F^{-1}\psi)) = \\ = f\left(F\left((2\pi)^{n/2} (F\phi F^{-1}\psi)\right)\right) & \\ = f\left(F\left((2\pi)^{n/2} (F^{-1} FF\phi F^{-1}\psi)\right)\right) &= f(F(F^{-1}(FF\phi * \psi))) \\ f(FF\phi * \psi) &= f(\psi * FF\phi). \end{aligned} \quad (13.19)$$

The last line follows from the following.

$$\begin{aligned} \int FF\phi(x - y) \psi(y) dy &= \int F\phi(x - y) F\psi(y) dy = \int F\psi(x - y) F\phi(y) dy \\ &= \int \psi(x - y) FF\phi(y) dy. \end{aligned}$$

From 13.19 and 13.18, since ψ was arbitrary,

$$(2\pi)^{n/2} F(F^{-1} \phi F^{-1} f) = (2\pi)^{n/2} F^{-1}(F\phi Ff) \equiv f * \phi$$

which shows 13.17. ■

13.3 Exercises

1. For $f \in L^1(\mathbb{R}^n)$, show that if $F^{-1}f \in L^1$ or $Ff \in L^1$, then f equals a continuous bounded function a.e.
2. Suppose $f, g \in L^1(\mathbb{R})$ and $Ff = Fg$. Show $f = g$ a.e.
3. Show that if $f \in L^1(\mathbb{R}^n)$, then $\lim_{|x| \rightarrow \infty} Ff(x) = 0$.
4. \uparrow Suppose $f * f = f$ or $f * f = 0$ and $f \in L^1(\mathbb{R})$. Show $f = 0$.
5. For this problem define $\int_a^\infty f(t) dt \equiv \lim_{r \rightarrow \infty} \int_a^r f(t) dt$. Note this coincides with the Lebesgue integral when $f \in L^1(a, \infty)$. Show

- (a) $\int_0^\infty \frac{\sin(u)}{u} du = \frac{\pi}{2}$
- (b) $\lim_{r \rightarrow \infty} \int_\delta^\infty \frac{\sin(ru)}{u} du = 0$ whenever $\delta > 0$.
- (c) If $f \in L^1(\mathbb{R})$, then $\lim_{r \rightarrow \infty} \int_{\mathbb{R}} \sin(ru) f(u) du = 0$.

Hint: For the first two, use $\frac{1}{u} = \int_0^\infty e^{-ut} dt$ and then, using this, apply Fubini's theorem to $\int_0^R \sin u \int_{\mathbb{R}} e^{-ut} dt du$. For the last part, first establish it for $f \in C_c^\infty(\mathbb{R})$ and then use the density of this set in $L^1(\mathbb{R})$ to obtain the result. This is sometimes called the Riemann Lebesgue lemma.

6. \uparrow Suppose that $g \in L^1(\mathbb{R})$ and that at some $x > 0$, g is locally Holder continuous from the right and from the left. This means

$$\lim_{r \rightarrow 0+} g(x+r) \equiv g(x+), \quad \lim_{r \rightarrow 0+} g(x-r) \equiv g(x-)$$

exists and there exist constants $K, \delta > 0$ and $r \in (0, 1]$ such that for $|x-y| < \delta$,

$$|g(x+) - g(y)| < K|x-y|^r, \quad |g(x-) - g(y)| < K|x-y|^r$$

for $y > x$ and $y < x$ respectively. Show that under these conditions,

$$\lim_{r \rightarrow \infty} \frac{2}{\pi} \int_0^\infty \frac{\sin(ur)}{u} \left(\frac{g(x-u) + g(x+u)}{2} \right) du = \frac{g(x+) + g(x-)}{2}.$$

7. Let $g \in L^1(\mathbb{R})$ and suppose g is locally Holder continuous from the right and from the left at x . Show that then

$$\lim_{R \rightarrow \infty} \frac{1}{2\pi} \int_{-R}^R e^{ixt} \int_{-\infty}^\infty e^{-ity} g(y) dy dt = \frac{g(x+) + g(x-)}{2}.$$

This is very interesting. If $g \in L^2(\mathbb{R})$, this shows $F^{-1}(Fg)(x) = \frac{g(x+) + g(x-)}{2}$, the midpoint of the jump in g at the point, x . In particular, if $g \in \mathcal{G}$, $F^{-1}(Fg) = g$. **Hint:** Show the left side of the above equation reduces to

$$\frac{2}{\pi} \int_0^\infty \frac{\sin(ur)}{u} \left(\frac{g(x-u) + g(x+u)}{2} \right) du$$

and then use Problem 6 to obtain the result.

8. \uparrow A measurable function g defined on $(0, \infty)$ has exponential growth if $|g(t)| \leq Ce^{\eta t}$ for some η . For $\operatorname{Re}(s) > \eta$, define the Laplace Transform: $Lg(s) \equiv \int_0^\infty e^{-su} g(u) du$. Assume that g has exponential growth as above and is Holder continuous from the right and from the left at t . Pick $\gamma > \eta$. Show that

$$\lim_{R \rightarrow \infty} \frac{1}{2\pi} \int_{-R}^R e^{\gamma y} e^{iyt} Lg(\gamma + iy) dy = \frac{g(t+) + g(t-)}{2}.$$

This formula is sometimes written in the form $\frac{1}{2\pi i} \int_{\gamma-i\infty}^{\gamma+i\infty} e^{st} Lg(s) ds$ and is called the complex inversion integral for Laplace transforms. It can be used to find inverse Laplace transforms. **Hint:** Plug in the formula for the Laplace transform and then massage to get it in the form of the preceding problem.

9. Suppose $f \in \mathcal{G}$. Show $F(f_{x_j})(t) = it_j Ff(t)$.
10. Let $f \in \mathcal{G}$ and let k be a positive integer. $\|f\|_{k,2} \equiv (\|f\|_2^2 + \sum_{|\alpha| \leq k} \|D^\alpha f\|_2^2)^{1/2}$. One could also define $\|f\|'_{k,2} \equiv (\int_{\mathbb{R}^n} |Ff(\mathbf{x})|^2 (1 + |\mathbf{x}|^2)^k dx)^{1/2}$. Show both $\|\cdot\|_{k,2}$ and $\|\cdot\|'_{k,2}$ are norms on \mathcal{G} and that they are equivalent. These are Sobolev space norms. For which values of k does the second norm make sense? How about the first norm?
11. \uparrow Define $H^k(\mathbb{R}^n)$, $k \geq 0$ by $f \in L^2(\mathbb{R}^n)$ such that

$$\left(\int |Ff(\mathbf{x})|^2 (1 + |\mathbf{x}|^2)^k dx \right)^{1/2} < \infty, \quad \|f\|'_{k,2} \equiv \left(\int |Ff(\mathbf{x})|^2 (1 + |\mathbf{x}|^2)^k dx \right)^{1/2}.$$

Show $H^k(\mathbb{R}^n)$ is a Banach space, and that if k is a positive integer, $H^k(\mathbb{R}^n)$ will be the set of all $f \in L^2(\mathbb{R}^n)$ such that there exists $\{u_j\} \subseteq \mathcal{G}$ with $\|u_j - f\|_2 \rightarrow 0$ and $\{u_j\}$ is a Cauchy sequence in $\|\cdot\|_{k,2}$ of Problem 10. This is one way to define Sobolev Spaces. **Hint:** One way to do the second part of this is to define a new measure μ by $\mu(E) \equiv \int_E (1 + |\mathbf{x}|^2)^k dx$. Then show μ is a Borel measure which is inner and outer regular and show there exists $\{g_m\}$ such that $g_m \in \mathcal{G}$ and $g_m \rightarrow Ff$ in $L^2(\mu)$. Thus $g_m = Ff_m$, $f_m \in \mathcal{G}$ because F maps \mathcal{G} onto \mathcal{G} . Then by Problem 10, $\{f_m\}$ is Cauchy in the norm $\|\cdot\|_{k,2}$.

12. \uparrow If $2k > n$, show that if $f \in H^k(\mathbb{R}^n)$, then f equals a bounded continuous function a.e. **Hint:** Show that for k this large, $Ff \in L^1(\mathbb{R}^n)$, and then use Problem 1. To do this, write

$$|Ff(\mathbf{x})| = |Ff(\mathbf{x})| (1 + |\mathbf{x}|^2)^{\frac{k}{2}} (1 + |\mathbf{x}|^2)^{-\frac{k}{2}},$$

So $\int |Ff(\mathbf{x})| dx = \int |Ff(\mathbf{x})| (1 + |\mathbf{x}|^2)^{\frac{k}{2}} (1 + |\mathbf{x}|^2)^{-\frac{k}{2}} dx$. Use the Cauchy Schwarz inequality. This is an example of a Sobolev imbedding Theorem.

13. Let $u \in \mathcal{G}$. Then $Fu \in \mathcal{G}$ and so, in particular, it makes sense to form the integral,

$$\int_{\mathbb{R}} Fu(\mathbf{x}', x_n) dx_n$$

where $(\mathbf{x}', x_n) = \mathbf{x} \in \mathbb{R}^n$. For $u \in \mathcal{G}$, define $\gamma u(\mathbf{x}') \equiv u(\mathbf{x}', 0)$. Find a constant such that $F(\gamma u)(\mathbf{x}')$ equals this constant times the above integral. **Hint:** By the dominated convergence theorem

$$\int_{\mathbb{R}} Fu(\mathbf{x}', x_n) dx_n = \lim_{\varepsilon \rightarrow 0} \int_{\mathbb{R}} e^{-(\varepsilon x_n)^2} Fu(\mathbf{x}', x_n) dx_n.$$

Now use the definition of the Fourier transform and Fubini's theorem as required in order to obtain the desired relationship.

14. Let $h(x) = \left(\int_0^x e^{-t^2} dt \right)^2 + \left(\int_0^1 \frac{e^{-x^2(1+t^2)}}{1+t^2} dt \right)$. Show that $h'(x) = 0$ and $h(0) = \pi/4$.

Then let $x \rightarrow \infty$ to conclude that $\int_0^\infty e^{-t^2} dt = \sqrt{\pi}/2$. Show that $\int_{-\infty}^\infty e^{-t^2} dt = \sqrt{\pi}$ and that $\int_{-\infty}^\infty e^{-ct^2} dt = \frac{\sqrt{\pi}}{\sqrt{c}}$.

15. Recall that for f a function, $f_{\mathbf{y}}(\mathbf{x}) = f(\mathbf{x} - \mathbf{y})$. Find a relationship between $Ff_{\mathbf{y}}(\mathbf{t})$ and $Ff(\mathbf{t})$ given that $f \in L^1(\mathbb{R}^n)$.
16. For $f \in L^1(\mathbb{R}^n)$, simplify $Ff(\mathbf{t} + \mathbf{y})$.
17. For $f \in L^1(\mathbb{R}^n)$ and c a nonzero real number, show $Ff(c\mathbf{t}) = Fg(\mathbf{t})$ where $g(\mathbf{x}) = f\left(\frac{\mathbf{x}}{c}\right)$.
18. Suppose that $f \in L^1(\mathbb{R})$ and that $\int |x| |f(x)| dx < \infty$. Find a way to use the Fourier transform of f to compute $\int xf(x) dx$.
19. Let (Ω, \mathcal{F}, P) be a probability space and let $X : \Omega \rightarrow \mathbb{R}^n$ be a random variable. This means $X^{-1}(\text{open set}) \in \mathcal{F}$. Define a measure λ_X on the Borel sets of \mathbb{R}^n as follows. For E a Borel set, $\lambda_X(E) \equiv P(X^{-1}(E))$. Explain why this is well defined. Next explain why λ_X can be considered a Radon probability measure by completion. Explain why $\lambda_X \in \mathcal{G}^*$ if $\lambda_X(\psi) \equiv \int_{\mathbb{R}^n} \psi d\lambda_X$ where \mathcal{G} is the collection of functions used to define the Fourier transform.
20. Using the above problem, the characteristic function of this measure (random variable) is $\phi_X(\mathbf{y}) \equiv \int_{\mathbb{R}^n} e^{i\mathbf{x} \cdot \mathbf{y}} d\lambda_X$. Show this always exists for any such random variable and is continuous. Next show that for two random variables X, Y , $\lambda_X = \lambda_Y$ if and only if $\phi_X(\mathbf{y}) = \phi_Y(\mathbf{y})$ for all \mathbf{y} . In other words, show the distribution measures are the same if and only if the characteristic functions are the same. A lot more can be concluded by looking at characteristic functions of this sort. The important thing about these characteristic functions is that they always exist, unlike moment generating functions.
21. Show that $C_c(\mathbb{R}^m)$, the continuous functions with compact support, with the norm given by $\|f\| \equiv \sup\{|f(\mathbf{y})| : \mathbf{y} \in \mathbb{R}^m\}$ is separable. **Hint:** You might note that this space is a subset of $C_0(\mathbb{R}^m)$ which is separable.
22. Show that if μ and ν are two Radon measures defined on σ algebras, \mathcal{S}_μ and \mathcal{S}_ν , of subsets of \mathbb{R}^n and if $\mu(V) = \nu(V)$ for all V open, then $\mu = \nu$ and $\mathcal{S}_\mu = \mathcal{S}_\nu$. **Hint:** Since the two measures agree on open sets, it follows that the two measures agree on every G_δ and F_σ set. By Proposition 11.1.2 on Page 315, if $E \in \mathcal{S}_\mu$, then there exists F, G such that $F \subseteq E \subseteq G$ and $\mu(G) = \nu(G) = \nu(F) = \mu(F)$ with F an F_σ set and G a G_δ set. Use completeness of the measures.

Chapter 14

Integration on Manifolds

Till now, integrals have mostly pertained to measurable subsets of \mathbb{R}^p and not something like a surface contained in a higher dimensional space. This is what is considered in this chapter. First is an abstract description of manifolds and then an interesting application of the representation theorem for positive linear functionals is used to give a measure on a manifold. This is the higher dimensional version of arc length for a smooth curve seen in calculus.

Definition 14.0.1 Let S be a nonempty set in a metric space (X, d) . ∂S is the set of points x , if any with the property that $B(x, r)$ contains points of S and points of $X \setminus S$ for each $r > 0$. The interior of S consists of the union of all open subsets of S .

Lemma 14.0.2 Let U be a nonempty open set in a metric space (X, d) . $\partial U = \bar{U} \setminus U$.

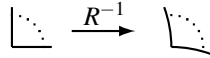
Proof: If $x \in \partial U$, then x can't be in U because some ball containing x is contained in U . However, it must be in \bar{U} because if not, some ball containing x would contain no points of \bar{U} since \bar{U} is closed.

If $x \in \bar{U} \setminus U$ then if some ball containing x fails to contain other points which are in U then that ball would show $x \notin \bar{U}$. Hence every ball containing x must contain points of U . However, x itself is not in U and so $x \in \partial U$. ■

14.1 Manifolds

Definition 14.1.1 An essential part of the definition of a manifold is the idea of a relatively open set defined next. Recall that a homeomorphism is a one to one, onto, continuous mapping from one metric space to another which has continuous inverse. A half space will be of the form $\{x : x_i \geq a_i\}$ or $\{x : x_i \leq a_i\}$.

Definition 14.1.2 Let X be a metric space and let $\Omega \subseteq X$. Then a set U is called a relatively open set or open in Ω if it is the intersection of an open set of X with Ω . Thus Ω is a metric space with respect to the distance $d(x, y)$ inherited from X and all considerations such as limit points, closures, limits, closed sets, open sets etc. in this metric space are taken with respect to this metric. Continuity is also defined in terms of this metric on Ω inherited from X . Ω is a p dimensional manifold with boundary if there is a locally finite cover $\{U_i\}$ (here it will be a finite cover) of sets open in Ω such that each U_i is homeomorphic to a set open in H where H is a half space or some finite intersection of such half spaces. Denote the open sets and homeomorphisms by (U_i, R_i) . The collection of these is called an atlas. Thus $R_i U_i$ is a set open in H_{R_i} where H_{R_i} is described above. Note that it could be a closed box. Then a point x is called a boundary point if and only if $R_i x$ is a boundary point of the interior of some H_{R_i} for some i .



I will be assuming that Ω is compact and so we can replace “locally finite” with finite in the above definition. First I need to verify that the idea of $\partial \Omega$ is well defined.

Lemma 14.1.3 $\partial \Omega$ is well defined in the sense that the statement that x is a boundary point does not depend on which chart is considered.

Proof: Suppose x is not a boundary point with respect to the chart (U, R) but is a boundary point with respect to (V, S) . Then $U \cap V$ is open in Ω so $Rx \in B \subseteq R(U \cap V)$ where $R(U \cap V)$ is open in H_R and B is an open ball contained in $R(U \cap V)$. But then, by Theorem 8.10.5, $S \circ R^{-1}(B)$ is open in \mathbb{R}^p and contains Sx so x is not a boundary point with respect to (V, S) after all. ■

Definition 14.1.4 Let $V \subseteq \mathbb{R}^q$. $C^k(\bar{V}; \mathbb{R}^p)$ is the set of functions which are restrictions to V of some function defined on \mathbb{R}^q which has k continuous derivatives which has values in \mathbb{R}^p . When $k = 0$, it means the restriction to V of continuous functions. A function is in $D(\bar{V}; \mathbb{R}^p)$ if it is the restriction to V of a differentiable function defined on \mathbb{R}^q . A Lipschitz function f is one which satisfies $\|f(x) - f(y)\| \leq K\|x - y\|$.

Thus, if $f \in C^k(\bar{V}; \mathbb{R}^q)$ or $D(\bar{V}; \mathbb{R}^p)$, we can consider it defined on \bar{V} and not just on V . This is the way one can generalize a one sided derivative of a function defined on a closed interval.

Lemma 14.1.5 Suppose A is a $m \times n$ matrix in which $m > n$ and A is one to one. Then $\|v\| \equiv |Av|$ is a norm on \mathbb{R}^n equivalent to the usual norm.

Proof: All the algebraic properties of the norm are obvious. If $\|v\| = 0$ then $|Av| = 0$ and since A is one to one, it follows $v = 0$ also. Now recall that all norms on \mathbb{R}^n are equivalent. ■

We have in mind, from now on that our manifold will be a compact subset of \mathbb{R}^q for some $q \geq p$.

Proposition 14.1.6 Suppose in the atlas for a manifold with boundary Ω it is also the case that each chart (U, R) has $R^{-1} \in C^1(\bar{R(U)})$ and $DR^{-1}(x)$ is one to one on $\bar{R(U)}$. Then for two charts (U, R) and (V, S) , it will be the case that $S \circ R^{-1} : R(U \cap V) \rightarrow S(V)$ will be also $C^1(\bar{R(U \cap V)})$.

Proof: Then

$$\begin{aligned} DR^{-1}(x)h + o(h) &= R^{-1}(x+h) - R^{-1}(x) \\ &= S^{-1}(S(R^{-1}(x+h))) - S^{-1}(S(R^{-1}(x))) \end{aligned} \quad (14.1)$$

$$\begin{aligned} &= DS^{-1}(S(R^{-1}(x)))(S(R^{-1}(x+h)) - S(R^{-1}(x))) \\ &\quad + o(S(R^{-1}(x+h)) - S(R^{-1}(x))) \end{aligned} \quad (14.2)$$

By continuity of R^{-1}, S , if h is small enough, which will always be assumed,

$$\begin{aligned} &|o(S(R^{-1}(x+h)) - S(R^{-1}(x)))| \\ &\leq \frac{\alpha}{2} |S(R^{-1}(x+h)) - S(R^{-1}(x))| \end{aligned}$$

where here there is $\alpha > 0$ such that

$$\begin{aligned} &|DS^{-1}(S(R^{-1}(x)))(S(R^{-1}(x+h)) - S(R^{-1}(x)))| \\ &\geq \alpha |S(R^{-1}(x+h)) - S(R^{-1}(x))| \end{aligned}$$

thanks to the assumption that $DS^{-1}(S(R^{-1}(x)))$ is one to one. Thus from 14.2

$$\frac{\alpha}{2} |(S(R^{-1}(x+h)) - S(R^{-1}(x)))| \leq |DR^{-1}(x)h + o(h)| \quad (14.3)$$

Now

$$\begin{aligned} & \frac{|o(S(R^{-1}(x+h)) - S(R^{-1}(x)))|}{|h|} \\ & \leq \frac{|o(S(R^{-1}(x+h)) - S(R^{-1}(x)))|}{|S(R^{-1}(x+h)) - S(R^{-1}(x))|} \frac{|S(R^{-1}(x+h)) - S(R^{-1}(x))|}{|h|} \end{aligned}$$

From 14.3, the second factor in the above is bounded. Now continuity of $S \circ R^{-1}$ implies that as $h \rightarrow 0$, the first factor also converges to 0. Thus

$$o(S(R^{-1}(x+h)) - S(R^{-1}(x))) = o(h)$$

Returning to 14.2,

$$DR^{-1}(x)h + o(h) = DS^{-1}(S(R^{-1}(x)))(S \circ R^{-1}(x+h) - S \circ R^{-1}(x))$$

Thus if $h = tv$,

$$\begin{aligned} & \lim_{t \rightarrow 0} DS^{-1}(S(R^{-1}(x))) \left(\frac{(S \circ R^{-1}(x+tv) - S \circ R^{-1}(x))}{t} \right) \\ & = DR^{-1}(x)v + \lim_{t \rightarrow 0} \frac{o(tv)}{t} = DR^{-1}(x)v \end{aligned}$$

By the above lemma, $\lim_{t \rightarrow 0} \frac{(S \circ R^{-1}(x+tv) - S \circ R^{-1}(x))}{t} = D_v(S \circ R^{-1})(x)$ exists. Also

$$DS^{-1}(S(R^{-1}(x)))D_v(S \circ R^{-1})(x) = DR^{-1}(x)v$$

Let $A(x) \equiv DS^{-1}(S(R^{-1}(x)))$. Then A^*A is invertible and $x \rightarrow A(x)$ is continuous. Then

$$\begin{aligned} A(x)^*A(x)D_v(S \circ R^{-1})(x) &= A(x)^*DR^{-1}(x)v \\ D_v(S \circ R^{-1})(x) &= (A(x)^*A(x))^{-1}A(x)^*DR^{-1}(x)v \end{aligned}$$

so $D_v(S \circ R^{-1})(x)$ is continuous. It follows from Theorem 7.6.1 that $S \circ R^{-1}$ is a function in $C^1(R(\overline{U \cap V}))$ because the Gateaux derivatives exist and are continuous. ■

Saying $DR^{-1}(x)$ is one to one is the analog of the situation in calculus with a smooth curve in which we assume the derivative is non zero and that the parametrization has continuous derivative.

I will be assuming that we can replace “locally finite” with finite in the above definition. This would happen, for example if Ω were compact, but this is not necessary. First I need to verify that the idea of $\partial\Omega$ is well defined.

Definition 14.1.7 A compact subset Ω of \mathbb{R}^q will be called a differentiable p -dimensional manifold with boundary if it is a C^0 manifold and also has some differentiable

structure about to be described. Ω is a differentiable manifold if $\mathbf{R}_j \circ \mathbf{R}_i^{-1}$ is differentiable on $\mathbf{R}_i(U_j \cap U_i)$. This is implied by the condition of Proposition 14.1.6. If, in addition to this, it has an atlas (U_i, \mathbf{R}_i) such that all partial derivatives are continuous and for all \mathbf{x}

$$\det(D\mathbf{R}_i^{-1}(\mathbf{R}_i(\mathbf{x})))^* (D\mathbf{R}_i^{-1}(\mathbf{R}_i(\mathbf{x}))) \neq 0$$

then it is called a smooth manifold. This condition is like the one for a smooth curve in calculus in which the derivative does not vanish. If, in addition “differentiable” is replaced with C^k meaning the first k derivatives exist and are continuous, then it will be a smooth C^k manifold with boundary.

Next is the concept of an oriented manifold. Orientation can be defined for general C^0 manifolds using the topological degree, but the reason for considering this, at least here, involves some sort of differentiability.

Definition 14.1.8 A differentiable manifold Ω with boundary is called orientable if there exists an atlas, $\{(U_r, \mathbf{R}_r)\}_{r=1}^m$, such that whenever $U_i \cap U_j \neq \emptyset$,

$$\det(D(\mathbf{R}_j \circ \mathbf{R}_i^{-1}))(\mathbf{u}) \geq 0 \text{ for all } \mathbf{u} \in \mathbf{R}_i(U_i \cap U_j) \quad (14.4)$$

An atlas satisfying 14.4 is called an oriented atlas. Also the following notation is often used with the convention that $\mathbf{v} = \mathbf{R}_i \circ \mathbf{R}_j^{-1}(\mathbf{u})$

$$\frac{\partial(v_1 \cdots v_p)}{\partial(u_1 \cdots u_p)} \equiv \det D(\mathbf{R}_i \circ \mathbf{R}_j^{-1})(\mathbf{u})$$

In this case, another atlas will be called an equivalent atlas (V_i, \mathbf{S}_i) if

$$\det(D(\mathbf{S}_j \circ \mathbf{R}_i^{-1}))(\mathbf{u}) \geq 0 \text{ for all } \mathbf{u} \in \mathbf{R}_i(U_i \cap V_j)$$

You can verify using the chain rule that this condition does indeed define an equivalence relation. Thus an oriented manifold would consist of a metric space along with an equivalence class of atlases. You could also define a piecewise smooth manifold as the union of finitely many smooth manifolds which have intersection only at boundary points.

Orientation is about the order in which the variables are listed or the way the positive coordinate axes point relative to each other. When you have an $n \times n$ matrix, you can always write its row reduced echelon form as a product of elementary matrices, some of which are permutation matrices or involve changing the direction by multiplying by a negative scalar, which also changes orientation the others having positive determinant. If there are an odd number of switches or multiplication by a negative scalar, you get the determinant is non-positive. If an even number, the determinant is non-negative. This is why we use the determinant to keep track of orientation in the above definition.

Example 14.1.9 Let $f: \mathbb{R}^{p+1} \rightarrow \mathbb{R}$ is C^1 and suppose and that $Df(\mathbf{x}) \neq \mathbf{0}$ for all \mathbf{x} contained in the set $\{\mathbf{x} : f(\mathbf{x}) = 0\}$. Then if $\{\mathbf{x} : f(\mathbf{x}) = 0\}$ is nonempty, it is a C^1 manifold thanks to an application of the implicit function theorem.

Note that this includes S^{p-1} , $\{\mathbf{x} \in \mathbb{R}^p : |\mathbf{x}| = 1\}$ and lots of other things like $x^4 + y^2 + z^4 = 1$ and so forth. The details are left as an exercise.

Recall from calculus how you can get pointy places in a space curve when the derivative of the parametrization is allowed to vanish. Here this would correspond to some $D\mathbf{R}_i^{-1}(\mathbf{u})$ not being one to one which is the same as having $D(\mathbf{R}_i \circ \mathbf{R}_i^{-1}(\mathbf{u}))$ having zero determinant.

In the above, it is not assumed that $D\mathbf{R}^{-1}$ is one to one. This can be used to include the concept of a higher dimensional version of a piecewise smooth curve. Suppose, for example you have $Q_1 \equiv [-1, 0] \times \prod_{i=2}^p [a_i, b_i]$, $Q_2 \equiv [0, 1] \times \prod_{i=2}^p [a_i, b_i]$ so there are two boxes joined along a common side. Let $\mathbf{R}_{1i}^{-1}, \mathbf{R}_{2j}^{-1}$ be as described above on these boxes and that \mathbf{R}_{1i}^{-1} and \mathbf{R}_{2j}^{-1} are continuous along the common face. We assume the union of $\mathbf{R}_{ri}^{-1}(U_{ri}), r = 1, 2$ is a smooth manifold so that $D\mathbf{R}_{ri}^{-1}$ exists on Q_r . Maybe $D\mathbf{R}_{1i}^{-1}, D\mathbf{R}_{2j}^{-1}$ are one to one on Q_1, Q_2 but on the common face, there is a difference in $D_1 \mathbf{R}_{1i}^{-1}, D_1 \mathbf{R}_{2j}^{-1}$ at a point on that face. Thus, if the restriction of \mathbf{R}_i^{-1} to Q_r is \mathbf{R}_{ri}^{-1} then \mathbf{R}_i^{-1} is not differentiable at points on this face. However, we could change the parametrization at the expense of allowing $D\mathbf{R}_{ri}^{-1}$ to equal zero on the common face which will result in \mathbf{R}_i^{-1} being differentiable. One simply replaces $\mathbf{x} \rightarrow \mathbf{R}_{ri}^{-1}(x_1, \dots, x_p)$ with $\mathbf{x} \rightarrow \mathbf{R}_{ri}^{-1}(x_1^3, x_2, \dots, x_p)$. This could be generalized to strings of boxes, successive pairs intersecting along a face thereby obtaining a higher dimensional notion of “piecewise smooth” as a case where the determinant of $D\mathbf{R}_i^{-1}$ is allowed to vanish. This is why it is useful in what follows to have a change of variables formula which does not require the non-vanishing of the determinant of the derivative of the transformation. This is the higher dimensional notion of pointy places occurring in space curves at points where the derivative vanishes. Note that the resulting union of the two smooth manifolds would end up being orientable if $\det(D(\mathbf{R}_{1j} \circ \mathbf{R}_{2i}^{-1}))(\mathbf{u}) > 0$ for all pertinent \mathbf{u} on the common face. Here we would take the partial derivative D_1 from the appropriate side in the chain rule. This is all very fussy but is mentioned to illustrate that in order to include piecewise smooth manifolds it suffices to only require that an atlas be differentiable. Thanks to Theorem 11.7.1 edges of a differentiable manifold can be ignored in the development of the area measure on a manifold if they result from some lower dimensional curve in \mathbb{R}^p or more generally a set of measure zero in \mathbb{R}^p . In this regard, see the rank theorem, Theorem 8.8.3 which identifies this as happening when $D\mathbf{R}_i^{-1}$ has smaller rank.

14.2 The Area Measure on a Manifold

Next the “surface measure” on a manifold is given. In what follows, the manifold will be a compact subset of \mathbb{R}^q . This has nothing to do with orientation. It will involve the following definition. To motivate this definition, recall the way you found the length of a curve in calculus where $t \in [a, b]$. It was $\int_a^b |\mathbf{r}'(t)| dt = \int_a^b \det(D\mathbf{r}(t)^* D\mathbf{r}(t))^{1/2} dt$ where $\mathbf{r}(t)$ is a parametrization for the curve. Think of $dl = \det(D\mathbf{r}(t)^* D\mathbf{r}(t))^{1/2} dt$ and you sum these to get the length.

Definition 14.2.1 Let (U_i, \mathbf{R}_i) be an atlas for a p dimensional differentiable manifold with boundary Ω . Also let $\{\psi_i\}_{i=1}^r$ be a partition of unity from Theorem 3.12.5 spt $\psi_i \subseteq U_i$. Then for $f \in C_c(\Omega)$, define

$$Lf \equiv \sum_{i=1}^r \int_{\mathbf{R}_i(U_i)} f(\mathbf{R}_i^{-1}(\mathbf{u})) \psi_i(\mathbf{R}_i^{-1}(\mathbf{u})) J_i(\mathbf{u}) d\mathbf{u}$$

Here du signifies $dm_p(u)$ and

$$J_i(u) \equiv (\det(DR_i^{-1}(u)^* DR_i^{-1}(u)))^{1/2}$$

I need to show that the same thing is obtained if another atlas and/or partition of unity is used.

Theorem 14.2.2 *The functional L is well defined in the sense that if another atlas is used, then for $f \in C_c(\Omega)$, the same value is obtained for Lf .*

Proof: Let the other atlas be $\{(V_j, S_j)\}_{j=1}^s$ where $v \in V_j$ and S_j has the same properties as the R_i . Then $(S_j \circ R_i^{-1})(u) = v$ so $R_i^{-1}(u) = S_j^{-1}(v)$ and so $R_i^{-1}(u) = S_j^{-1}((S_j \circ R_i^{-1})(u))$ implying $DR_i^{-1}(u) = DS_j^{-1}(v)D(S_j \circ R_i^{-1})(u)$. Therefore,

$$\begin{aligned} J_i(u) &= (\det(DR_i^{-1}(u)^* DR_i^{-1}(u)))^{1/2} \\ &= \left(\det \left(\overbrace{D(S_j \circ R_i^{-1})^*(u)}^{p \times p} \overbrace{DS_j^{-1}(v)^* DS_j^{-1}(v)}^{(p \times q)(q \times p)} \overbrace{D(S_j \circ R_i^{-1})(u)}^{p \times p} \right) \right)^{1/2} \\ &= \left[\det(D(S_j \circ R_i^{-1})^*(u)) \det(D(S_j \circ R_i^{-1})(u)) \right]^{1/2} J_j(v) \\ &= |\det(D(S_j \circ R_i^{-1})(u))| J_j(v) \end{aligned} \quad (14.5)$$

Similarly

$$J_j(v) = |\det(D(R_i \circ S_j^{-1})(v))| J_i(u). \quad (14.6)$$

Let \hat{L} go with this new atlas. Thus

$$\hat{L}(f) \equiv \sum_{j=1}^s \int_{S_j(V_j)} f(S_j^{-1}(v)) \eta_j(S_j^{-1}(v)) J_j(v) dv \quad (14.7)$$

where η_j is a partition of unity associated with the sets V_j as described above. Now letting ψ_i be the partition of unity for the U_i , $v = S_j \circ R_i^{-1}(u)$ for $u \in R_i(V_j \cap U_i)$.

$$\begin{aligned} & \int_{S_j(V_j)} f(S_j^{-1}(v)) \eta_j(S_j^{-1}(v)) J_j(v) dv \\ &= \sum_{i=1}^r \int_{S_j(V_j \cap U_i)} f(S_j^{-1}(v)) \psi_i(S_j^{-1}(v)) \eta_j(S_j^{-1}(v)) J_j(v) dv \end{aligned}$$

By Lemma 11.8.1, the assumptions of differentiability imply that the boundary points of Ω are always mapped to a set of measure zero so these can be neglected if desired. Now $S_j(V_j \cap U_i) = S_j \circ R_i^{-1}(R_i(V_j \cap U_i))$ and so using 14.6, the above expression equals

$$\begin{aligned} & \sum_{i=1}^r \int_{R_i(V_j \cap U_i)} f(R_i^{-1}(u)) \psi_i(R_i^{-1}(u)) \eta_j(R_i^{-1}(u)) \cdot \\ & \quad \left| \det(D(R_i \circ S_j^{-1})(v)) \right| J_i(u) |\det D(S_j \circ R_i^{-1})(u)| du \end{aligned}$$

Now $I = (\mathbf{R}_i \circ \mathbf{S}_j^{-1}) \circ (\mathbf{S}_j \circ \mathbf{R}_i^{-1})$ and so the chain rule implies that the product of the two Jacobians is 1. Hence 14.7 equals

$$\begin{aligned}
 & \sum_{j=1}^s \sum_{i=1}^r \int_{\mathbf{R}_i(V_j \cap U_i)} f(\mathbf{R}_i^{-1}(\mathbf{u})) \psi_i(\mathbf{R}_i^{-1}(\mathbf{u})) \eta_j(\mathbf{R}_i^{-1}(\mathbf{u})) J_i(\mathbf{u}) d\mathbf{u} \\
 &= \sum_{i=1}^r \sum_{j=1}^s \int_{\mathbf{R}_i(U_i)} f(\mathbf{R}_i^{-1}(\mathbf{u})) \psi_i(\mathbf{R}_i^{-1}(\mathbf{u})) \eta_j(\mathbf{R}_i^{-1}(\mathbf{u})) J_i(\mathbf{u}) d\mathbf{u} \\
 &= \sum_{i=1}^r \int_{\mathbf{R}_i(U_i)} f(\mathbf{R}_i^{-1}(\mathbf{u})) \psi_i(\mathbf{R}_i^{-1}(\mathbf{u})) \sum_{j=1}^s \eta_j(\mathbf{R}_i^{-1}(\mathbf{u})) J_i(\mathbf{u}) d\mathbf{u} \\
 &= \sum_{i=1}^r \int_{\mathbf{R}_i(U_i)} f(\mathbf{R}_i^{-1}(\mathbf{u})) \psi_i(\mathbf{R}_i^{-1}(\mathbf{u})) J_i(\mathbf{u}) d\mathbf{u} = L(f)
 \end{aligned}$$

Thus L is a well defined positive linear functional. ■

Definition 14.2.3 By the representation theorem for positive linear functionals, Theorem 11.2.2, there exists a complete Radon measure σ_p defined on the Borel sets of Ω such that $Lf = \int_{\Omega} f d\sigma_p$. Then σ_p is what is meant by the measure on the differentiable manifold Ω .

If O is an open set in Ω , what is $\sigma_p(O)$? Let $f_n \uparrow \mathcal{X}_O$ where f_n is continuous. Then by the monotone convergence theorem,

$$\begin{aligned}
 \sigma_p(O) &= \lim_{n \rightarrow \infty} L(f_n) = \lim_{n \rightarrow \infty} \sum_{i=1}^r \int_{\mathbf{R}_i(U_i)} f_n(\mathbf{R}_i^{-1}(\mathbf{u})) \psi_i(\mathbf{R}_i^{-1}(\mathbf{u})) J_i(\mathbf{u}) d\mathbf{u} \\
 &= \lim_{n \rightarrow \infty} \sum_{i=1}^r \int_{\mathbf{R}_i(U_i \cap O)} f_n(\mathbf{R}_i^{-1}(\mathbf{u})) \psi_i(\mathbf{R}_i^{-1}(\mathbf{u})) J_i(\mathbf{u}) d\mathbf{u} \\
 &= \sum_{i=1}^r \int_{\mathbf{R}_i(U_i \cap O)} \mathcal{X}_O(\mathbf{R}_i^{-1}(\mathbf{u})) \psi_i(\mathbf{R}_i^{-1}(\mathbf{u})) J_i(\mathbf{u}) d\mathbf{u}.
 \end{aligned}$$

If K is a compact subset of some U_i , then use Corollary 12.6.5 to obtain a partition of unity which has $\psi_i = 1$ on K so that all other ψ_j equal 0. Then

$$\int_{\Omega} \mathcal{X}_K d\sigma_p = \int_{\mathbf{R}_i(U_i)} \mathcal{X}_K(\mathbf{R}_i^{-1}(\mathbf{u})) J_i(\mathbf{u}) d\mathbf{u}$$

It then follows from regularity of the measure and the monotone convergence theorem that if E is any measurable set contained in U_i , you can replace K in the above with E . In general, this implies that for nonnegative measurable f , having support in U_i ,

$$\int_{\Omega} f d\sigma_p = \int_{\mathbf{R}_i(U_i)} f(\mathbf{R}_i^{-1}(\mathbf{u})) J_i(\mathbf{u}) d\mathbf{u}$$

Indeed, $\partial\Omega$ is a closed subset of Ω and so $\mathcal{X}_{\partial\Omega}$ is measurable. That part of the boundary contained in U_i would then involve a Lebesgue integral over a set of m_p measure zero. This shows the following proposition.

Proposition 14.2.4 Let Ω be a differentiable manifold as discussed above and let σ_p be the measure on the manifold defined above. Then $\sigma_p(\partial\Omega) = 0$.

Using Rademacher's theorem which is presented later, which says that every Lipschitz function defined on an open set is differentiable a.e. and Theorem 11.7.1 which says that Lipschitz functions map sets of measure zero to sets of measure zero and measurable sets to measurable sets, the following corollary is obtained using the same arguments.

Corollary 14.2.5 *Let Ω be a subset of a finite dimensional normed linear space be a Lipschitz manifold meaning that it is a C^0 manifold for which each atlas (U_i, \mathbf{R}_i) has \mathbf{R}_i Lipschitz on U_i and \mathbf{R}_i^{-1} is Lipschitz on $\mathbf{R}_i(U_i)$. Then there is a regular complete measure σ_p defined on a σ algebra \mathcal{F} of subsets of Ω which is finite on compact sets, includes the Borel sets, and satisfies*

$$\int_{\Omega} f d\sigma_p = \sum_{i=1}^r \int_{\mathbf{R}_i(U_i)} f(\mathbf{R}_i^{-1}(\mathbf{u})) \psi_i(\mathbf{R}_i^{-1}(\mathbf{u})) J_i(\mathbf{u}) d\mathbf{u}$$

for all $f \in C_c(\Omega)$ where here the partition of unity comes from Theorem 3.12.5 with respect to the distance coming from the norm.

The justification for using $J_i(\mathbf{u}) d\mathbf{u}$ in the definition of the area measure comes from geometric reasoning which has not been presented yet. This will be done in the chapter on Hausdorff measures. However, it was noted above that this is a generalization of a familiar example from calculus. It would also be possible to verify that it works from familiar definitions in calculus in the case of a two dimensional manifold. Also note that it suffices to assume only that $D\mathbf{R}_i^{-1}(\mathbf{u})$ exists for a.e. \mathbf{u} .

14.3 Divergence Theorem

The divergence theorem considered here will feature an open set in \mathbb{R}^p whose boundary has a particular form. For convenience, if $\mathbf{x} \in \mathbb{R}^p$, $\hat{\mathbf{x}}_i \equiv (x_1 \cdots x_{i-1} \ x_{i+1} \cdots x_p)^T$.

Definition 14.3.1 *Let $U \subseteq \mathbb{R}^p$ satisfy the following conditions. There exist open boxes, Q_1, \dots, Q_N , $Q_i = \prod_{j=1}^p (a_j^i, b_j^i)$ such that $\partial U \equiv \overline{U} \setminus U$ is contained in their union. Also, there exists an open set, Q_0 such that $Q_0 \subseteq \overline{Q_0} \subseteq U$ and $\overline{U} \subseteq Q_0 \cup Q_1 \cup \dots \cup Q_N$. Assume for each Q_i , there exists k and a function g_i such that $U \cap Q_i$ is of the form*

$$\left\{ \begin{array}{l} \mathbf{x} : (x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_p) \in \prod_{j=1}^{k-1} (a_j^i, b_j^i) \times \\ \prod_{j=k+1}^p (a_j^i, b_j^i) \text{ and } a_k^i < x_k < g_i(x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_p) \end{array} \right\} \quad (14.8)$$

or else of the form

$$\left\{ \begin{array}{l} \mathbf{x} : (x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_p) \in \prod_{j=1}^{k-1} (a_j^i, b_j^i) \times \\ \prod_{j=k+1}^p (a_j^i, b_j^i) \text{ and } g_i(x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_p) < x_k < b_k^i \end{array} \right\} \quad (14.9)$$

The function, g_i is differentiable and has a measurable partial derivatives on

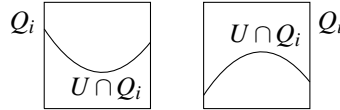
$$A_i \subseteq \prod_{j=1}^{k-1} (a_j^i, b_j^i) \times \prod_{j=k+1}^p (a_j^i, b_j^i) \equiv \hat{Q}_k$$

where

$$m_{p-1} \left(\prod_{j=1}^{k-1} (a_j^i, b_j^i) \times \prod_{j=k+1}^p (a_j^i, b_j^i) \setminus A_i \right) = 0.$$

and we assume there is a constant C such that for all i and j , $\left| \frac{\partial g_i}{\partial x_j} \right| \leq C$ off A_i and that in each variable, g can be recovered from integrating an appropriate partial derivative. That is, each g_k is absolutely continuous in each variable.

To illustrate the above here is a picture.



Recall from calculus that if $z - g(\hat{x}) = 0$ then to get a normal vector to the level surface, it will be \pm the gradient.

Lemma 14.3.2 Let $\alpha_1, \dots, \alpha_p$ be real numbers and let $A(\alpha_1, \dots, \alpha_p)$ be the matrix which has $1 + \alpha_i^2$ in the i^{th} slot and $\alpha_i \alpha_j$ in the ij^{th} slot when $i \neq j$. Then $\det A = 1 + \sum_{i=1}^p \alpha_i^2$.

Proof of the claim: The matrix, $A(\alpha_1, \dots, \alpha_p)$ is of the form

$$A(\alpha_1, \dots, \alpha_p) = \begin{pmatrix} 1 + \alpha_1^2 & \alpha_1 \alpha_2 & \cdots & \alpha_1 \alpha_p \\ \alpha_1 \alpha_2 & 1 + \alpha_2^2 & & \alpha_2 \alpha_p \\ \vdots & & \ddots & \vdots \\ \alpha_1 \alpha_p & \alpha_2 \alpha_p & \cdots & 1 + \alpha_p^2 \end{pmatrix}$$

Now consider the product of a matrix and its transpose, $B^T B$ below.

$$\begin{pmatrix} 1 & 0 & \cdots & 0 & \alpha_1 \\ 0 & 1 & & 0 & \alpha_2 \\ \vdots & & \ddots & & \vdots \\ 0 & & & 1 & \alpha_p \\ -\alpha_1 & -\alpha_2 & \cdots & -\alpha_p & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & \cdots & 0 & -\alpha_1 \\ 0 & 1 & & 0 & -\alpha_2 \\ \vdots & & \ddots & & \vdots \\ 0 & & & 1 & -\alpha_p \\ \alpha_1 & \alpha_2 & \cdots & \alpha_p & 1 \end{pmatrix} \quad (14.10)$$

This product equals a matrix of the form $\begin{pmatrix} A(\alpha_1, \dots, \alpha_p) & \mathbf{0} \\ \mathbf{0} & 1 + \sum_{i=1}^p \alpha_i^2 \end{pmatrix}$. Therefore, $(1 + \sum_{i=1}^p \alpha_i^2) \det(A(\alpha_1, \dots, \alpha_p)) = \det(B)^2 = \det(B^T)^2$. However, using row operations,

$$\det B^T = \det \begin{pmatrix} 1 & 0 & \cdots & 0 & \alpha_1 \\ 0 & 1 & & 0 & \alpha_2 \\ \vdots & & \ddots & & \vdots \\ 0 & & & 1 & \alpha_p \\ 0 & 0 & \cdots & 0 & 1 + \sum_{i=1}^p \alpha_i^2 \end{pmatrix} = 1 + \sum_{i=1}^p \alpha_i^2$$

and therefore,

$$\left(1 + \sum_{i=1}^p \alpha_i^2 \right) \det(A(\alpha_1, \dots, \alpha_p)) = \left(1 + \sum_{i=1}^p \alpha_i^2 \right)^2$$

which shows $\det(A(\alpha_1, \dots, \alpha_p)) = (1 + \sum_{i=1}^p \alpha_i^2)$. ■

Now consider the case of σ on ∂U . The maps will be of the form

$$\hat{x} \in Q_k \rightarrow (x_1 \quad \cdots \quad x_{i-1} \quad g(\hat{x}_i) \quad x_{i+1} \quad \cdots \quad x_p)^T = h(\hat{x}_i)$$

I need to describe $\det(Dh(\hat{x}_i)^* Dh(\hat{x}_i))^{1/2} \equiv J(\hat{x})$.

Consider an example sufficient to see what happens in general in which $p = 3$ and $i = 2$. Then in this case, $J(\hat{x})$ will be the square root of the determinant of

$$\begin{pmatrix} 1 & g_{x_1} & 0 \\ 0 & g_{x_3} & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ g_{x_1} & g_{x_3} \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} g_{x_1}^2 + 1 & g_{x_1} g_{x_3} \\ g_{x_1} g_{x_3} & g_{x_3}^2 + 1 \end{pmatrix}.$$

One can verify that this is just a special case in which $Dh(\hat{x}_i)^* Dh(\hat{x}_i)$ will be of the form considered in Lemma 14.3.2. Thus by this lemma, $J(x) = \sqrt{1 + \sum_{k \neq i} g_{x_k}^2}$.

Then if $U \cap Q$ is of the form in 14.8 or in 14.9 one can identify the unit exterior normal to the surface either on the top or the bottom of $U \cap Q$ from beginning calculus. These are respectively

$$\mathbf{n} = \frac{(-g_{x_1} \quad \cdots \quad -g_{x_{p-1}} \quad \cdots \quad 1)^T}{\sqrt{1 + \sum_{k=1}^{p-1} g_{x_k}^2}}, \quad \frac{(g_{x_1} \quad \cdots \quad g_{x_{p-1}} \quad \cdots \quad -1)^T}{\sqrt{1 + \sum_{k=1}^{p-1} g_{x_k}^2}}$$

The first pointing up away from U and the second pointing down away from U .

If you simply assume g_k is differentiable, there is no problem in Definition 14.3.1. One can show with Rademacher's theorem that it suffices to assume these functions are Lipschitz continuous.

In the following proof, I will regard $f(x_1, x_2, \dots, x_p)$ as a function of the listed variables.

Definition 14.3.3 Let $\mathbf{F} \in C^1(\bar{U}; \mathbb{R}^p)$ and the rectangular coordinates are denoted as $\mathbf{x} = (x_1, \dots, x_p)$. Then the divergence of \mathbf{F} written as $\text{div}(\mathbf{F})$ is defined as $\sum_i \frac{\partial F_i}{\partial x_i} \equiv \sum_i F_{i,i}$. It is also written as $\nabla \cdot \mathbf{F}$.

Theorem 14.3.4 Let U be a bounded open set in \mathbb{R}^p satisfying the conditions of Definition 14.3.1 and let $\mathbf{F} \in C^1(\bar{U}; \mathbb{R}^p)$. Then

$$\int_U \sum_{i=1}^p F_{i,i}(\mathbf{x}) dm_p = \int_{\partial U} \mathbf{F} \cdot \mathbf{n} d\sigma_{p-1}$$

where \mathbf{n} is the unit exterior normal to U just described.

Proof: Let $\text{spt}(\psi_i)$ be a compact subset of Q_i and $\sum_{i=0}^N \psi_i = 1$ on \bar{U} and each ψ_i is infinitely differentiable. This partition of unity exists by Lemma 12.6.4. There is an explicit description of the unit outer normal for each point of the boundary of U described above in either of the two cases described in Definition 14.3.1 and illustrated in the above picture. Then

$$\begin{aligned} \int_U \sum_i F_{i,i}(\mathbf{x}) dm_p &= \int_U \sum_i \sum_{k=0}^N (\psi_k F)_{i,i}(\mathbf{x}) dm_p = \sum_i \sum_{k=0}^N \int_U (\psi_k F)_{i,i}(\mathbf{x}) dm_p \\ &= \sum_{k=0}^N \int_{Q_k} \sum_i (\psi_k F)_{i,i}(\mathbf{x}) dm_p \end{aligned} \quad (14.11)$$

Now consider one of the terms in the above. For the sake of simplicity assume $k = p$ so that the special direction corresponds to x_p . Also, I will assume that the function $g(\hat{x})$ is on the top, so it is like the left picture in the above. A similar argument works if $g(\hat{x})$ were on the bottom. Either way we can specify a unit exterior normal a.e. I will omit the subscript on g_k , Q_k , and ψ_k .

Case that $i < p$: Pick $i < p$. Letting \hat{Q} be (x_1, \dots, x_{p-1}) where $x \in Q$, For any i ,

$$\int_Q (\psi_k F)_{i,i} dm_p = \int_{\hat{Q}} \int_{-\infty}^{g(x_1, \dots, x_{p-1})} (\psi_k F_i)_i dx_p d\hat{x} = \int_{\hat{Q}} \int_{-\infty}^0 D_i(\psi F_i)(\hat{x}, y + g(\hat{x})) dy d\hat{x} \quad (14.12)$$

Now for $i < p$, that in the integrand is not $\frac{\partial}{\partial x_i}(\psi F_i)(\hat{x}, y + g(\hat{x}))$. Indeed, by the chain rule,

$$\frac{\partial}{\partial x_i}(\psi F_i)(\hat{x}, y + g(\hat{x})) = D_i(\psi F_i)(\hat{x}, y + g(\hat{x})) + D_p(\psi F_i)(\hat{x}, y + g(\hat{x})) \frac{\partial g(\hat{x})}{\partial x_i}$$

Since $\text{spt}(\psi) \subseteq Q$, it follows that 14.12 reduces to

$$\begin{aligned} \int_{-\infty}^0 \int_{\hat{Q}} \frac{\partial}{\partial x_i}(\psi F_i)(\hat{x}, y + g(\hat{x})) d\hat{x} dy - \int_{\hat{Q}} \int_{-\infty}^0 D_p(\psi F_i)(\hat{x}, y + g(\hat{x})) \frac{\partial g(\hat{x})}{\partial x_i} dy d\hat{x} \\ = 0 - \int_{\hat{Q}} (\psi F_i)(\hat{x}, g(\hat{x})) d\hat{x} \end{aligned}$$

Case that $i = p$: In this case, 14.12 becomes $\int_{\hat{Q}} (\psi F_p)(\hat{x}, g(\hat{x})) d\hat{x}$. Recall how it was just shown that the unit normal is $\frac{(-g_{x_1}, \dots, -g_{x_{p-1}}, 1)}{\sqrt{\sum_{i=1}^{p-1} g_{x_k}^2 + 1}}$ and $d\sigma = \sqrt{\sum_{i=1}^{p-1} g_{x_k}^2 + 1} dm_{p-1}$. Then the above reduces to $\int_{\partial(Q \cap U)} (\psi F) \cdot n d\sigma$. The same result will hold for all the Q_i . The sign changes if in the situation of 14.9. As to Q_0 , $\int_{Q_0} \sum_i (\psi_0 F)_{i,i}(x) dm_p = 0$ because $\text{spt}(\psi_0) \subseteq Q_0$. Returning to 14.11, it follows that

$$\begin{aligned} \int_U \sum_i F_{i,i}(x) dm_p &= \sum_{k=0}^N \int_{Q_k} \sum_i (\psi_k F)_{i,i}(x) dm_p = \sum_{k=0}^N \int_{Q_k} \sum_i (\psi_k F)_{i,i}(x) dm_p \\ &= \sum_{k=1}^N \int_{\partial(Q_k \cap U)} (\psi_k F) \cdot n d\sigma = \sum_{k=1}^N \int_{\partial U} (\psi_k F) \cdot n d\sigma \\ &= \int_{\partial U} \left(\sum_{k=0}^N \psi_k \right) F \cdot n d\sigma = \int_{\partial U} F \cdot n d\sigma \quad \blacksquare \end{aligned}$$

Definition 14.3.5 The expression $\sum_{i=1}^p F_{i,i}(x)$ is called $\text{div}(F)$. It is defined above in terms of the coordinates with respect to a fixed orthonormal basis (e_1, \dots, e_p) . However, it does not depend on such a particular choice for coordinates.

If you had some other orthonormal basis (v_1, \dots, v_p) and if (y_1, \dots, y_p) are the coordinates of a point z with respect to this other orthonormal system, then there is an orthogonal matrix Q such that $y = Qx$ for y the coordinate vector for the new basis and x the coordinate vector for the old basis. Then

$$J_i(x) \equiv (\det(DR_i^{-1}(x)^* DR_i^{-1}(x)))^{1/2} = \left(\det \left((DR_i^{-1}(y)Q)^* DR_i^{-1}(y)Q \right) \right)^{1/2}$$

$$= (\det(Q^* DR_i^{-1}(\mathbf{y})^* DR_i^{-1}(\mathbf{y})Q))^{1/2} = (\det(Q^* DR_i^{-1}(\mathbf{y})^* DR_i^{-1}(\mathbf{y})Q))^{1/2} = J_i(\mathbf{y})$$

so the two definitions of $d\sigma$ will be the same with either set of coordinates.

List the \mathbf{v}_i in the order which will give $\det(Q) = 1$. That is to say, the two bases have the same orientation. The insistence that $\det Q = 1$ will ensure that the unit normal vectors defined as above will point away from U . Thus we could take the divergence with respect to coordinates of any orthonormal basis having the same orientation. Note that for a.e. geometric point \mathbf{z}

$$\operatorname{div}(\mathbf{F})(\mathbf{z}) = \lim_{r \rightarrow 0} \frac{1}{m_p(B(\mathbf{z}, r))} \int_{B(\mathbf{z}, r)} \operatorname{div}(\mathbf{F}) dm_p = \lim_{r \rightarrow 0} \frac{1}{m_p(B(\mathbf{z}, r))} \int_{\partial B(\mathbf{z}, r)} \mathbf{F} \cdot \mathbf{n} d\sigma_{p-1}$$

the first equal sign from the fundamental theorem of calculus and the last expression on the right being independent of the choice of basis. This implies that we could have generalized the kind of region to be one for which the little rectangles are allowed to be slanted. Creases and pointy places in the manifold can result from places where some $J_i(\mathbf{x}) = 0$, due to some DR_i^{-1} not being one to one, but this will not matter because in the definition of the surface measure this will be a set of measure zero on the manifold. The change of variables formula which was so important in the above argument is unaffected by these creases.

Globally the region could be quite complicated. As an example in two dimensions, it might look like this:



Corollary 14.3.6 *If the divergence is computed with respect to \mathbf{y} where $\mathbf{y} = Q\mathbf{x}$ for Q orthogonal with determinant 1, and each box used in the argument of Theorem 14.3.4 is taken with respect to such a new basis $(\mathbf{v}_1, \dots, \mathbf{v}_p)$, then one still obtains $\int_U \operatorname{div}(\mathbf{F}) dm_p = \int_{\partial U} \mathbf{F} \cdot \mathbf{n} d\sigma_{p-1}$.*

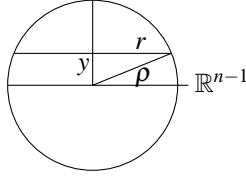
14.4 Volumes of Balls in \mathbb{R}^p

This short section will give an explicit description of surface area given in Section 11.11.

Recall, $B(\mathbf{x}, r)$ denotes the set of all $\mathbf{y} \in \mathbb{R}^p$ such that $|\mathbf{y} - \mathbf{x}| < r$. By the change of variables formula for multiple integrals or simple geometric reasoning, all balls of radius r have the same volume. Furthermore, simple reasoning or change of variables formula will show that the volume of the ball of radius r equals $\alpha_p r^p$ where α_p will denote the volume of the unit ball in \mathbb{R}^p . With the divergence theorem, it is now easy to give a simple relationship between the surface area of the ball of radius r and the volume. Let $d\alpha_{p-1}$ be the area measure above. By the divergence theorem, $\int_{B(\mathbf{0}, r)} \operatorname{div}(\mathbf{x}) dx = \int_{\partial B(\mathbf{0}, r)} \mathbf{x} \cdot \frac{\mathbf{x}}{|\mathbf{x}|} d\alpha_{p-1}$ because the unit outward normal on $\partial B(\mathbf{0}, r)$ is $\frac{\mathbf{x}}{|\mathbf{x}|}$. Therefore, $p\alpha_p r^p = r\alpha_{p-1}(\partial B(\mathbf{0}, r))$ and so $\alpha_{p-1}(\partial B(\mathbf{0}, r)) = p\alpha_p r^{p-1}$.

Let ω_p denote the area of the sphere $S^{p-1} = \{\mathbf{x} \in \mathbb{R}^p : |\mathbf{x}| = 1\}$. I just showed that $\omega_p = p\alpha_p$.

I want to find α_p now.



Taking slices at height y as shown and using that these slices have $p-1$ dimensional area equal to $\alpha_{p-1}r^{p-1}$, it follows $\alpha_p\rho^p = 2\int_0^\rho \alpha_{p-1}(\rho^2 - y^2)^{(p-1)/2} dy$ since the r at a given y is $\sqrt{\rho^2 - y^2}$. In the integral, change variables, letting $y = \rho \cos \theta$. Then $\alpha_p\rho^p = 2\rho^p\alpha_{p-1}\int_0^{\pi/2} \sin^p(\theta) d\theta$. It follows that

$$\alpha_p = 2\alpha_{p-1} \int_0^{\pi/2} \sin^p(\theta) d\theta. \quad (14.13)$$

From this we find a formula for α_p .

First note that $\Gamma(\frac{1}{2}) = \int_0^\infty e^{-t} t^{-1/2} dt = \int_0^\infty e^{-u^2} u^{-1} 2u du = 2 \int_0^\infty e^{-u^2} = \sqrt{\pi}$ from elementary calculus using polar coordinates and change of variables.

Theorem 14.4.1 $\alpha_p = \frac{\pi^{p/2}}{\Gamma(\frac{p}{2}+1)}$ where Γ denotes the gamma function, defined for $\alpha > 0$ by $\Gamma(\alpha) \equiv \int_0^\infty e^{-t} t^{\alpha-1} dt$.

Proof: Let $p = 1$ first. Then $\alpha_1 = \pi = \frac{\pi^{1/2}}{\Gamma(\frac{1}{2}+1)}$ because $\Gamma(\alpha+1) = \alpha\Gamma(\alpha)$ so the right side is $\frac{\pi^{1/2}}{\frac{1}{2}\Gamma(\frac{1}{2})} = 2$ which is indeed the one dimensional area of the unit ball in one dimension. Similarly it is true for $p = 2, 3$. Assume true for $p \geq 3$. Then using 14.13 and induction,

$$\alpha_{p+1} = 2 \frac{\alpha_p}{\Gamma(\frac{p}{2}+1)} \int_0^{\pi/2} \sin^{p+1}(\theta) d\theta$$

Using an integration by parts, this equals $2 \frac{\pi^{p/2}}{\Gamma(\frac{p}{2}+1)} \frac{p}{p+1} \int_0^{\pi/2} \sin^{p-1}(\theta) d\theta$. By 14.13 and induction this is

$$\begin{aligned} \frac{\pi^{p/2}}{\Gamma(\frac{p}{2}+1)} \frac{p}{p+1} \frac{\alpha_{p-1}}{\alpha_{p-2}} &= \frac{\pi^{p/2}}{\Gamma(\frac{p}{2}+1)} \frac{p}{p+1} \frac{\frac{\pi^{(p-1)/2}}{\Gamma(\frac{p-1}{2}+1)}}{\frac{\pi^{(p-2)/2}}{\Gamma(\frac{p-2}{2}+1)}} = \frac{2\pi^{(p+1)/2} \Gamma(\frac{p}{2})}{\Gamma(\frac{p}{2}+1) \Gamma(\frac{p-1}{2}+1)} \frac{p/2}{p+1} \\ &= \frac{2\pi^{(p+1)/2} \Gamma(\frac{p}{2}+1)}{\Gamma(\frac{p}{2}+1) \Gamma(\frac{p-1}{2}+1)} \frac{1}{p+1} = \frac{\pi^{(p+1)/2}}{\Gamma(\frac{p+1}{2})} \frac{1}{\frac{p+1}{2}} = \frac{\pi^{(p+1)/2}}{\Gamma(\frac{p+1}{2}+1)} \blacksquare \end{aligned}$$

14.5 Exercises

1. A random vector \mathbf{X} , with values in \mathbb{R}^p has a multivariate normal distribution written as $\mathbf{X} \sim N_p(\mathbf{m}, \Sigma)$ if for all Borel $E \subseteq \mathbb{R}^p$,

$$\lambda_{\mathbf{X}}(E) = \int_{\mathbb{R}^p} \mathcal{X}_E(\mathbf{x}) \frac{1}{(2\pi)^{p/2} \det(\Sigma)^{1/2}} e^{\frac{-1}{2}(\mathbf{x}-\mathbf{m})^* \Sigma^{-1}(\mathbf{x}-\mathbf{m})} d\mathbf{m}_p$$

Here Σ is a positive definite symmetric matrix. Recall that $\lambda_{\mathbf{X}}(E) \equiv P(\mathbf{X} \in E)$. Using the change of variables formula, show that $\lambda_{\mathbf{X}}$ defined above is a probability measure. One thing you must show is that

$$\int_{\mathbb{R}^p} \frac{1}{(2\pi)^{p/2} \det(\Sigma)^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\mathbf{m})^* \Sigma^{-1}(\mathbf{x}-\mathbf{m})} d\mathbf{m}_p = 1$$

Hint: To do this, you might use the fact from linear algebra that $\Sigma = Q^* D Q$ where D is a diagonal matrix and Q is an orthogonal matrix. Thus $\Sigma^{-1} = Q^* D^{-1} Q$. Maybe you could first let $\mathbf{y} = D^{-1/2} Q(\mathbf{x} - \mathbf{m})$ and change the variables. Note that the change of variables formula works fine when the open sets are all of \mathbb{R}^p . You don't need to confine your attention to finite open sets which would be the case with Riemann integrals which are only defined on bounded sets.

2. Consider the surface $z = x^2$ for $(x, y) \in (0, 1) \times (0, 1)$. Find the area of this surface. **Hint:** You can make do with just one chart in this case. Let $\mathbf{R}^{-1}(x, y) = (x, y, x^2)^T$, $(x, y) \in (0, 1) \times (0, 1)$. Then

$$D\mathbf{R}^{-1} = \begin{pmatrix} 1 & 0 & 2x \\ 0 & 1 & 0 \end{pmatrix}^T$$

$$\text{It follows that } D\mathbf{R}^{-1*} D\mathbf{R}^{-1} = \begin{pmatrix} 4x^2 + 1 & 0 \\ 0 & 1 \end{pmatrix}.$$

3. A parametrization for most of the sphere of radius $a > 0$ in three dimensions is

$$\begin{aligned} x &= a \sin(\phi) \cos(\theta) \\ y &= a \sin(\phi) \sin(\theta) \\ z &= a \cos(\phi) \end{aligned}$$

where we will let $\phi \in (0, \pi)$, $\theta \in (0, 2\pi)$ so there is just one chart involved. As mentioned earlier, this includes all of the sphere except for the line of longitude corresponding to $\theta = 0$. Find a formula for the area of this sphere. Again, we are making do with a single chart.

4. Let V be such that the divergence theorem holds. Show that $\int_V \nabla \cdot (u \nabla v) dV = \int_{\partial V} u \frac{\partial v}{\partial \mathbf{n}} dA$ where \mathbf{n} is the exterior normal and $\frac{\partial v}{\partial \mathbf{n}}$ denotes the directional derivative of v in the direction \mathbf{n} . Remember the directional derivative.

$$\lim_{t \rightarrow 0} \frac{v(\mathbf{x} + t\mathbf{n}) - v(\mathbf{x})}{t} = \lim_{t \rightarrow 0} \frac{Dv(\mathbf{x})(t\mathbf{n}) + o(t)}{t} = Dv(\mathbf{x})(\mathbf{n}) = \nabla v(\mathbf{x}) \cdot \mathbf{n}$$

5. To prove the divergence theorem, it was shown first that the spacial partial derivative in the volume integral could be exchanged for multiplication by an appropriate component of the exterior normal. This problem starts with the divergence theorem and goes the other direction. Assuming the divergence theorem, holds for a region V , show that $\int_{\partial V} \mathbf{n} u dA = \int_V \nabla u dV$. Note this implies $\int_V \frac{\partial u}{\partial x} dV = \int_{\partial V} n_1 u dA$.
6. Fick's law for diffusion states the flux of a diffusing species, \mathbf{J} is proportional to the gradient of the concentration c . Write this law getting the sign right for the constant of proportionality and derive an equation similar to the heat equation for the concentration c . Typically, c is the concentration of some sort of pollutant or a chemical.

7. Sometimes people consider diffusion in materials which are not homogeneous. This means that $\mathbf{J} = -K\nabla c$ where K is a 3×3 matrix and c is called the concentration. Thus in terms of components, $J_i = -\sum_j K_{ij} \frac{\partial c}{\partial x_j}$. Here c is the concentration which means the amount of pollutant or whatever is diffusing in a volume is obtained by integrating c over the volume. Derive a formula for a nonhomogeneous model of diffusion based on the above.
8. Let V be such that the divergence theorem holds. Show that

$$\int_V (v\nabla^2 u - u\nabla^2 v) dV = \int_{\partial V} \left(v \frac{\partial u}{\partial n} - u \frac{\partial v}{\partial n} \right) dA$$

where \mathbf{n} is the exterior normal and $\frac{\partial u}{\partial n}$ is defined in Problem 4. Here $\nabla^2 u \equiv \sum_i u_{,x_i x_i}$.

9. Let V be a ball and suppose $\nabla^2 u = f$ in V while $u = g$ on ∂V . Show that there is at most one solution to this boundary value problem which is C^2 in V and continuous on V with its boundary. **Hint:** You might consider $w = u - v$ where u and v are solutions to the problem. Then use the result of Problem 4 and the identity $w\nabla^2 w = \nabla \cdot (w\nabla w) - \nabla w \cdot \nabla w$ to conclude $\nabla w = 0$. Then show this implies w must be a constant by considering $h(t) = w(t\mathbf{x} + (1-t)\mathbf{y})$ and showing h is a constant.
10. Show that $\int_{\partial V} \nabla \times \mathbf{v} \cdot \mathbf{n} dA = 0$ where V is a region for which the divergence theorem holds and \mathbf{v} is a C^2 vector field.
11. Let $\mathbf{F}(x, y, z) = (x, y, z)$ be a vector field in \mathbb{R}^3 and let V be a three dimensional shape and let $\mathbf{n} = (n_1, n_2, n_3)$. Show that $\int_{\partial V} (xn_1 + yn_2 + zn_3) dA = 3 \times \text{volume of } V$.
12. Let $\mathbf{F} = x\mathbf{i} + y\mathbf{j} + z\mathbf{k}$ and let V denote the tetrahedron formed by the planes, $x = 0$, $y = 0$, $z = 0$, and $\frac{1}{3}x + \frac{1}{3}y + \frac{1}{3}z = 1$. Verify the divergence theorem for this example.
13. Suppose $f : U \rightarrow \mathbb{R}$ is continuous where U is some open set and for all $B \subseteq U$ where B is a ball, $\int_B f(\mathbf{x}) dV = 0$. Show that this implies $f(\mathbf{x}) = 0$ for all $\mathbf{x} \in U$.
14. Let U denote the box centered at $(0, 0, 0)$ with sides parallel to the coordinate planes which has width 4, length 2 and height 3. Find the flux integral $\int_{\partial U} \mathbf{F} \cdot \mathbf{n} dS$ where $\mathbf{F} = (x + 3, 2y, 3z)$. **Hint:** If you like, you might want to use the divergence theorem.
15. Find the flux out of the cylinder whose base is $x^2 + y^2 \leq 1$ which has height 2 of the vector field $\mathbf{F} = (xy, zy, z^2 + x)$.
16. Find the flux out of the ball of radius 4 centered at $\mathbf{0}$ of the vector field $\mathbf{F} = (x, zy, z + x)$.
17. In one dimension, the heat equation is of the form $u_t = \alpha u_{xx}$. Show that $u(x, t) = e^{-\alpha n^2 t} \sin(nx)$ satisfies the heat equation

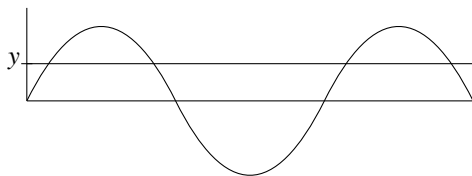
Chapter 15

Degree Theory

This chapter is on the Brouwer degree, a very useful concept with numerous and important applications. The degree can be used to prove some difficult theorems in topology such as the Brouwer fixed point theorem, the Jordan separation theorem, and the invariance of domain theorem. A couple of these big theorems have been presented earlier, but when you have degree theory, they get much easier. Degree theory is also used in bifurcation theory and many other areas in which it is an essential tool. The degree will be developed for \mathbb{R}^p in this book. When this is understood, it is not too difficult to extend to versions of the degree which hold in Banach space. There is more on degree theory in the book by Deimling [12] and much of the presentation here follows this reference. Another more recent book which is really good is [15]. This is a whole book on degree theory.

The original reference for the approach given here, based on analysis, is [27] and dates from 1959. The degree was developed earlier by Brouwer and others using different methods. The more classical approach based on simplices and approximations with these is in [29]. I have given an approach based on singular homology as an appendix in [38].

To give you an idea what the degree is about, consider a real valued C^1 function defined on an interval I , and let $y \in f(I)$ be such that $f'(x) \neq 0$ for all $x \in f^{-1}(y)$. In this case the degree is the sum of the signs of $f'(x)$ for $x \in f^{-1}(y)$, written as $d(f, I, y)$.



In the above picture, $d(f, I, y)$ is 0 because there are two places where the sign is 1 and two where it is -1 .

The amazing thing about this is the number you obtain in this simple manner is a specialization of something which is defined for continuous functions and which has nothing to do with differentiability. The reason one can extend the above simple idea to continuous functions is is an integral expression for the degree which is insensitive to homotopy. It is very similar to the winding number of complex analysis. The difference between the two is that with the degree, the integral which ties it all together is taken over the open set while the winding number is taken over the boundary, although proofs of it in the case of the winding number sometimes involve Green's theorem which involves an integral over the open set. I think these analogies are better seen in the other presentation in [38].

In this chapter Ω will refer to a bounded open set.

Definition 15.0.1 For Ω a bounded open set, denote by $C(\overline{\Omega})$ the set of functions which are restrictions of functions in $C_c(\mathbb{R}^p)$, equivalently $C(\mathbb{R}^p)$ to $\overline{\Omega}$ and by $C^m(\overline{\Omega})$, $m \leq \infty$ the space of restrictions of functions in $C_c^m(\mathbb{R}^p)$, equivalently $C^m(\mathbb{R}^p)$ to $\overline{\Omega}$. If $f \in C(\overline{\Omega})$ the symbol f will also be used to denote a function defined on \mathbb{R}^p equalling f on $\overline{\Omega}$ when convenient. The subscript c indicates that the functions have compact support. The norm in $C(\overline{\Omega})$ is defined as follows.

$$\|f\|_{\infty, \overline{\Omega}} = \|f\|_{\infty} \equiv \sup \{|f(x)| : x \in \overline{\Omega}\}.$$

If the functions take values in \mathbb{R}^p write $C^m(\bar{\Omega}; \mathbb{R}^p)$ or $C(\bar{\Omega}; \mathbb{R}^p)$ for these functions if there is no differentiability assumed. The norm on $C(\bar{\Omega}; \mathbb{R}^p)$ is defined in the same way as above,

$$\|f\|_{\infty, \bar{\Omega}} = \|f\|_{\infty} \equiv \sup \{|f(x)| : x \in \bar{\Omega}\}.$$

If $m = \infty$, the notation means that there are infinitely many derivatives. Also, $C(\Omega; \mathbb{R}^p)$ consists of functions which are continuous on Ω that have values in \mathbb{R}^p and $C^m(\Omega; \mathbb{R}^p)$ denotes the functions which have m continuous derivatives defined on Ω . Also let \mathcal{P} consist of functions $f(x)$ such that $f_k(x)$ is a polynomial, meaning an element of the algebra of functions generated by $\{1, x_1, \dots, x_p\}$. Thus a typical polynomial is of the form $\sum_{i_1 \dots i_p} a(i_1 \dots i_p) x^{i_1} \dots x^{i_p}$ where the i_j are nonnegative integers and $a(i_1 \dots i_p)$ is a real number.

Some of the theorems are simpler if you base them on the Weierstrass approximation theorem.

Note that, by applying the Tietze extension theorem to the components of the function, one can always extend a function continuous on $\bar{\Omega}$ to all of \mathbb{R}^p so there is no loss of generality in simply regarding functions continuous on $\bar{\Omega}$ as restrictions of functions continuous on \mathbb{R}^p . Next is the idea of a regular value.

Definition 15.0.2 For W an open set in \mathbb{R}^p and $g \in C^1(W; \mathbb{R}^p)$, y is called a regular value of g if whenever $x \in g^{-1}(y)$, $\det(Dg(x)) \neq 0$. Note that if $g^{-1}(y) = \emptyset$, it follows that y is a regular value from this definition. That is, y is a regular value if and only if

$$y \notin g(\{x \in W : \det Dg(x) = 0\})$$

Denote by S_g the set of singular values of g , those y such that $\det(Dg(x)) = 0$ for some $x \in g^{-1}(y)$.

Also, $\partial\Omega$ will often be referred to. It is those points with the property that every open set (or open ball) containing the point contains points not in Ω and points in Ω . Then the following simple lemma will be used frequently.

Lemma 15.0.3 Define ∂U to be those points x with the property that for every $r > 0$, $B(x, r)$ contains points of U and points of U^c . Then for U an open set, $\partial U = \bar{U} \setminus U$. Let C be a closed subset of \mathbb{R}^p and let \mathcal{K} denote the set of components of $\mathbb{R}^p \setminus C$. Then if K is one of these components, it is open and $\partial K \subseteq C$.

Proof: Let $x \in \bar{U} \setminus U$. If $B(x, r)$ contains no points of U , then $x \notin \bar{U}$. If $B(x, r)$ contains no points of U^c , then $x \in U$ and so $x \notin \bar{U} \setminus U$. Therefore, $\bar{U} \setminus U \subseteq \partial U$. Now let $x \in \partial U$. If $x \in U$, then since U is open there is a ball containing x which is contained in U contrary to $x \in \partial U$. Therefore, $x \notin U$. If x is not a limit point of U , then some ball containing x contains no points of U contrary to $x \in \partial U$. Therefore, $x \in \bar{U} \setminus U$ which shows the two sets are equal.

Why is K open for K a component of $\mathbb{R}^p \setminus C$? This follows from Theorem 3.11.12 and results from open balls being connected. Thus if $k \in K$, letting $B(k, r) \subseteq C^c$, it follows $K \cup B(k, r)$ is connected and contained in C^c and therefore is contained in K because K is maximal with respect to being connected and contained in C^c .

Now for K a component of $\mathbb{R}^p \setminus C$, why is $\partial K \subseteq C$? Let $x \in \partial K$. If $x \notin C$, then $x \in K_1$, some component of $\mathbb{R}^p \setminus C$. If $K_1 \neq K$ then x cannot be a limit point of K and so it cannot

be in ∂K . Therefore, $K = K_1$ but this also is a contradiction because if $x \in \partial K$ then $x \notin K$ thanks to the first part that $\partial U = \overline{U} \setminus U$. ■

Note that for an open set $U \subseteq \mathbb{R}^p$, and $h : \overline{U} \rightarrow \mathbb{R}^p$, $\text{dist}(h(\partial U), y) \geq \text{dist}(h(\overline{U}), y)$ because $\overline{U} \supseteq \partial U$.

The following lemma will be nice to keep in mind.

Lemma 15.0.4 $f \in C(\overline{\Omega} \times [a, b]; \mathbb{R}^p)$ if and only if

$$t \rightarrow f(\cdot, t) \in C([a, b]; C(\overline{\Omega}; \mathbb{R}^p))$$

Also

$$\|f\|_{\infty, \overline{\Omega} \times [a, b]} = \max_{t \in [a, b]} \left(\|f(\cdot, t)\|_{\infty, \overline{\Omega}} \right)$$

Proof: \Rightarrow By uniform continuity, if $\varepsilon > 0$ there is $\delta > 0$ such that if $|t - s| < \delta$, then for all $x \in \overline{\Omega}$, $\|f(x, t) - f(x, s)\| < \frac{\varepsilon}{2}$. It follows that

$$\|f(\cdot, t) - f(\cdot, s)\|_{\infty} \leq \frac{\varepsilon}{2} < \varepsilon$$

\Leftarrow Say $(x_n, t_n) \rightarrow (x, t)$. Does it follow that $f(x_n, t_n) \rightarrow f(x, t)$?

$$\begin{aligned} \|f(x_n, t_n) - f(x, t)\| &\leq \|f(x_n, t_n) - f(x_n, t)\| + \|f(x_n, t) - f(x, t)\| \\ &\leq \|f(\cdot, t_n) - f(\cdot, t)\|_{\infty} + \|f(x_n, t) - f(x, t)\| \end{aligned}$$

both terms converge to 0, the first because f is continuous into $C(\overline{\Omega}; \mathbb{R}^p)$ and the second because $x \rightarrow f(x, t)$ is continuous.

The claim about the norms is next. Let (x, t) be such that $\|f\|_{\infty, \overline{\Omega} \times [a, b]} < \|f(x, t)\| + \varepsilon$. Then

$$\|f\|_{\infty, \overline{\Omega} \times [a, b]} < \|f(x, t)\| + \varepsilon \leq \max_{t \in [a, b]} \left(\|f(\cdot, t)\|_{\infty, \overline{\Omega}} \right) + \varepsilon$$

and so $\|f\|_{\infty, \overline{\Omega} \times [a, b]} \leq \max_{t \in [a, b]} \left(\|f(\cdot, t)\|_{\infty, \overline{\Omega}} \right)$ because ε is arbitrary. However, the same argument works in the other direction. There exists t such that

$$\|f(\cdot, t)\|_{\infty, \overline{\Omega}} = \max_{t \in [a, b]} \left(\|f(\cdot, t)\|_{\infty, \overline{\Omega}} \right)$$

by compactness of the interval. Then by compactness of $\overline{\Omega}$, there is x such that

$$\|f(\cdot, t)\|_{\infty, \overline{\Omega}} = \|f(x, t)\| \leq \|f\|_{\infty, \overline{\Omega} \times [a, b]}$$

and so the two norms are the same. ■

15.1 Sard's Lemma and Approximation

First are easy assertions about approximation of continuous functions with smooth ones.

The following is the Weierstrass approximation theorem. It is Corollary 5.6.3 presented earlier.

Corollary 15.1.1 *If $f \in C([a, b]; X)$ where X is a normed linear space, then there exists a sequence of polynomials which converge uniformly to f on $[a, b]$. The polynomials are of the form*

$$\sum_{k=0}^m p_k(t) f\left(l\left(\frac{k}{m}\right)\right) \quad (15.1)$$

where l is a linear one to one and onto map from $[0, 1]$ to $[a, b]$ and $p_0(a) = 1$ but $p_k(a) = 0$ if $k \neq 0$, $p_m(b) = 1$ but $p_k(b) = 0$ for $k \neq m$.

Applying the Weierstrass approximation theorem, Theorem 5.8.7 or Theorem 5.10.5 to the components of a vector valued function yields the following Theorem.

Theorem 15.1.2 *If $f \in C(\overline{\Omega}; \mathbb{R}^p)$ for Ω a bounded subset of \mathbb{R}^p , then for any $\varepsilon > 0$, there exists $g \in C^\infty(\overline{\Omega}; \mathbb{R}^p)$ such that $\|g - f\|_{\infty, \overline{\Omega}} < \varepsilon$.*

Recall Sard's lemma, shown earlier. It is Lemma 11.8.3. I am stating it here for convenience.

Lemma 15.1.3 (Sard) *Let Ω be an open set in \mathbb{R}^p and let $h : \Omega \rightarrow \mathbb{R}^p$ be differentiable. Let*

$$S \equiv \{x \in \Omega : \det Dh(x) = 0\}.$$

Then $m_p(h(S)) = 0$.

First note that if $y \notin g(\Omega)$, then $y \notin g(\{x \in \Omega : \det Dg(x) = 0\})$ so it is a regular value.

Observe that any uncountable set in \mathbb{R}^p has a limit point. To see this, tile \mathbb{R}^p with countably many congruent boxes. One of them has uncountably many points. Now subdivide this into 2^p congruent boxes. One has uncountably many points. Continue subdividing this way to obtain a limit point as the unique point in the intersection of a nested sequence of compact sets whose diameters converge to 0.

Lemma 15.1.4 *Let $g \in C^\infty(\mathbb{R}^p; \mathbb{R}^p)$ and let $\{y_i\}_{i=1}^\infty$ be points of \mathbb{R}^p and let $\eta > 0$. Then there exists e with $\|e\| < \eta$ and $y_i + e$ is a regular value for g for all i .*

Proof: Let $S = \{x \in \mathbb{R}^p : \det Dg(x) = 0\}$. By Sard's lemma, $g(S)$ has measure zero. Let $N \equiv \bigcup_{i=1}^\infty (g(S) - y_i)$. Thus N has measure 0. Pick $e \in B(0, \eta) \setminus N$. Then for each i , $y_i + e \notin g(S)$. ■

Next we approximate f with a smooth function g such that each y_i is a regular value of g .

Lemma 15.1.5 *Let $f \in C(\overline{\Omega}; \mathbb{R}^p)$, Ω a bounded open set, and let $\{y_i\}_{i=1}^\infty$ be points not in $f(\partial\Omega)$ and let $\delta > 0$. Then there exists $g \in C^\infty(\overline{\Omega}; \mathbb{R}^p)$ such that $\|g - f\|_{\infty, \overline{\Omega}} < \delta$ and y_i is a regular value for g for each i . That is, if $g(x) = y_i$, then $Dg(x)^{-1}$ exists. Also, if $\delta < \text{dist}(f(\partial\Omega), y)$ for some y a regular value of $g \in C^\infty(\overline{\Omega}; \mathbb{R}^p)$, then $g^{-1}(y)$ is a finite set of points in Ω . Also, if y is a regular value of $g \in C^\infty(\mathbb{R}^p, \mathbb{R}^p)$, then $g^{-1}(y)$ is countable.*

Proof: Pick $\tilde{g} \in C^\infty(\bar{\Omega}; \mathbb{R}^p)$, $\|\tilde{g} - f\|_{\infty, \bar{\Omega}} < \delta$. From Lemma 15.1.4, $y_i + e$ is a regular value for \tilde{g} for each i where e can be chosen as small as desired. Let $g = \tilde{g} - e$ where e is so small that also $\|g - f\|_{\infty, \bar{\Omega}} < \delta$. Thus y_i is a regular value of g for all i . (same as $y_i + e$ regular value of \tilde{g}). This shows the first part.

It remains to verify the last claims. Since $\|g - f\|_{\infty, \bar{\Omega}} < \delta$, if $x \in \partial\Omega$, then

$$\|g(x) - y\| \geq \|f(x) - y\| - \|f(x) - g(x)\| \geq \text{dist}(f(\partial\Omega), y) - \delta > \delta - \delta = 0$$

and so $y \notin g(\partial\Omega)$, so if $g(x) = y$, then $x \in \Omega$. Thus $g^{-1}(y)$ is a compact subset of Ω and so for each $x \in g^{-1}(y)$ there is a ball containing x , B_x contained in Ω such that there is at most one point in $g^{-1}(y) \cap B_x$ this by the inverse function theorem. Finitely many of these balls cover $g^{-1}(y)$ so this set must be finite and at each point, the determinant of the derivative of g is nonzero. For y a regular value, $g^{-1}(y)$ is countable since otherwise, there would be a limit point $x \in g^{-1}(y)$ and g would fail to be one to one near x contradicting the inverse function theorem. ■

Now with this, here is a definition of the degree.

Definition 15.1.6 Let Ω be a bounded open set in \mathbb{R}^p and let $f : \bar{\Omega} \rightarrow \mathbb{R}^p$ be continuous. Let $y \notin f(\partial\Omega)$. Then the degree is defined as follows: Let g be infinitely differentiable,

$$\|f - g\|_{\infty, \bar{\Omega}} < \delta \equiv \text{dist}(f(\partial\Omega), y),$$

and y is a regular value of g . Then $y \notin g(\partial\Omega)$ and we define

$$d(f, \Omega, y) \equiv \sum \{ \text{sgn}(\det(Dg(x))) : x \in g^{-1}(y), x \in \Omega \}$$

where the sum is finite by Lemma 15.1.5, defined to equal 0 if $g^{-1}(y)$ is empty.

Note that if g is such an approximation of f then if $x \in \partial\Omega$ and $t \in [0, 1]$,

$$\begin{aligned} |tg(x) + (1-t)f(x) - y| &\geq |f(x) - y| - t\|g - f\|_{\infty} \\ &> \text{dist}(f(\partial\Omega), y) - \text{dist}(f(\partial\Omega), y) = 0 \end{aligned}$$

Thus $tg + (1-t)f$ maps no point of $\partial\Omega$ to y . In particular, g maps no point of $\partial\Omega$ to y .

Lemma 15.1.7 The above sum in the definition makes sense for a single g and, assuming this definition of $d(f, \Omega, y)$ is well defined, then it would follow that if $y \notin f(\Omega)$, then $d(f, \Omega, y) = 0$.

Proof: As just noted, if $\|f - g\|_{\infty, \bar{\Omega}} < \text{dist}(f(\partial\Omega), y)$ then $y \notin g(\partial\Omega)$. In fact

$$y \notin (tg(x) + (1-t)f(x))(\partial\Omega)$$

for any $t \in [0, 1]$. Thus the sum is a finite sum and makes sense by Lemma 15.1.5. What if $y \notin f(\Omega)$? In this case, assuming the definition is well defined, you could pick g such that y is a regular value for g and also $\|f - g\|_{\infty, \bar{\Omega}} < \text{dist}(f(\bar{\Omega}), y)$ so the above definition would say that $d(f, \Omega, y) = 0$ because there would be no terms in the sum. ■

We really need to verify that this definition is well defined, not dependent on which g is chosen. This involves the use of an integral.

Next is an identity. It was Lemma 7.11.2 on Page 201.

Lemma 15.1.8 Let $g : \Omega \rightarrow \mathbb{R}^p$ be C^2 where Ω is an open subset of \mathbb{R}^p . Then

$$\sum_{j=1}^p \operatorname{cof}(Dg)_{ij,j} = 0,$$

where here $(Dg)_{ij} \equiv g_{i,j} \equiv \frac{\partial g_i}{\partial x_j}$. Also, $\operatorname{cof}(Dg)_{ij} = \frac{\partial \det(Dg)}{\partial g_{i,j}}$.

Next is an integral representation of $\sum \{\operatorname{sgn}(\det(Dg(x))) : x \in g^{-1}(y)\}$ but first is a little lemma about disjoint sets.

Lemma 15.1.9 Let K be a compact set and C a closed set in \mathbb{R}^p such that $K \cap C = \emptyset$. Then

$$\operatorname{dist}(K, C) \equiv \inf \{\|k - c\| : k \in K, c \in C\} > 0.$$

Proof: Let $d \equiv \inf \{\|k - c\| : k \in K, c \in C\}$. Let $\{k_i\}, \{c_i\}$ be such that

$$d + \frac{1}{i} > \|k_i - c_i\|.$$

Since K is compact, there is a subsequence still denoted by $\{k_i\}$ such that $k_i \rightarrow k \in K$. Then also

$$\|c_i - c_m\| \leq \|c_i - k_i\| + \|k_i - k_m\| + \|c_m - k_m\|$$

If $d = 0$, then as $m, i \rightarrow \infty$ it follows $\|c_i - c_m\| \rightarrow 0$ and so $\{c_i\}$ is a Cauchy sequence which must converge to some $c \in C$. But then $\|c - k\| = \lim_{i \rightarrow \infty} \|c_i - k_i\| = 0$ and so $c = k \in C \cap K$, a contradiction to these sets being disjoint. ■

In particular the distance between a point and a closed set is always positive if the point is not in the closed set. Of course this is obvious even without the above lemma.

Definition 15.1.10 Let $g \in C^\infty(\overline{\Omega}; \mathbb{R}^p)$ where Ω is a bounded open set. Also let ϕ_ε be a mollifier.

$$\phi_\varepsilon \in C_c^\infty(B(0, \varepsilon)), \phi_\varepsilon \geq 0, \int \phi_\varepsilon dx = 1.$$

The idea is that ε will converge to 0 to get suitable approximations.

First, here is a technical lemma which will be used to identify the degree with an integral.

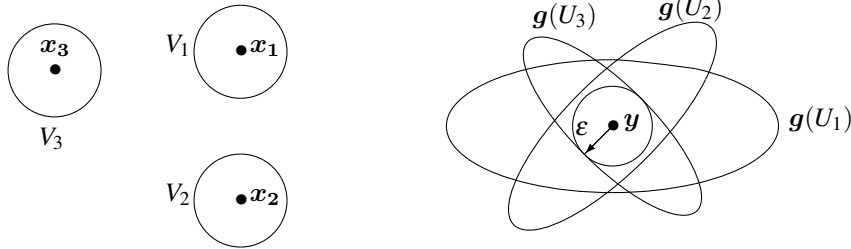
Lemma 15.1.11 Let $y \notin g(\partial\Omega)$ for $g \in C^\infty(\overline{\Omega}; \mathbb{R}^p)$. Also suppose y is a regular value of g . Then for all positive ε small enough,

$$\int_{\Omega} \phi_\varepsilon(g(x) - y) \det Dg(x) dx = \sum \{\operatorname{sgn}(\det Dg(x)) : x \in g^{-1}(y)\}$$

Proof: First note that the sum is finite from Lemma 15.1.5. It only remains to verify the equation. If $y \notin g(\Omega)$, then for $\varepsilon < \operatorname{dist}(g(\overline{\Omega}), y)$, $\phi_\varepsilon(g(x) - y) = 0$ for all $x \in \Omega$ so both sides equal 0.

I need to show the left side of this equation is constant for ε small enough and equals the right side. By what was just shown, there are finitely many points, $\{x_i\}_{i=1}^m = g^{-1}(y)$. By the inverse function theorem, there exist disjoint open sets U_i with $x_i \in U_i$, such that g is one

to one on U_i with $\det(Dg(x))$ having constant sign on U_i and $g(U_i)$ is an open set containing y . Then let ε be small enough that $B(y, \varepsilon) \subseteq \cap_{i=1}^m g(U_i)$. Also, $y \notin g(\overline{\Omega} \setminus (\cup_{i=1}^n U_i))$, a compact set. Let ε be still smaller, if necessary, so that $B(y, \varepsilon) \cap g(\overline{\Omega} \setminus (\cup_{i=1}^n U_i)) = \emptyset$ and let $V_i \equiv g^{-1}(B(y, \varepsilon)) \cap U_i$.



Therefore, for any ε this small,

$$\int_{\Omega} \phi_{\varepsilon}(g(x) - y) \det Dg(x) dx = \sum_{i=1}^m \int_{V_i} \phi_{\varepsilon}(g(x) - y) \det Dg(x) dx$$

The reason for this is as follows. The integrand on the left is nonzero only if $g(x) - y \in B(0, \varepsilon)$ which occurs only if $g(x) \in B(y, \varepsilon)$ which is the same as $x \in g^{-1}(B(y, \varepsilon))$. Therefore, the integrand is nonzero only if x is contained in exactly one of the disjoint sets, V_i . Now using the change of variables theorem, ($z = g(x) - y, g^{-1}(y + z) = x$.)

$$= \sum_{i=1}^m \int_{g(V_i) - y} \phi_{\varepsilon}(z) \det Dg(g^{-1}(y + z)) |\det Dg^{-1}(y + z)| dz \quad (15.2)$$

By the chain rule, $I = Dg(g^{-1}(y + z)) Dg^{-1}(y + z)$ and so in the above for a single V_i ,

$$\begin{aligned} & \det Dg(g^{-1}(y + z)) |\det Dg^{-1}(y + z)| \\ &= \operatorname{sgn}(\det Dg(g^{-1}(y + z))) |\det Dg(g^{-1}(y + z))| |\det Dg^{-1}(y + z)| \\ &= \operatorname{sgn}(\det Dg(g^{-1}(y + z))) = \operatorname{sgn}(\det Dg(x)) = \operatorname{sgn}(\det Dg(x_i)). \end{aligned}$$

Therefore, 15.2 reduces to

$$\begin{aligned} & \sum_{i=1}^m \operatorname{sgn}(\det Dg(x_i)) \int_{g(V_i) - y} \phi_{\varepsilon}(z) dz = \\ & \sum_{i=1}^m \operatorname{sgn}(\det Dg(x_i)) \int_{B(0, \varepsilon)} \phi_{\varepsilon}(z) dz = \sum_{i=1}^m \operatorname{sgn}(\det Dg(x_i)). \end{aligned}$$

In case $g^{-1}(y) = \emptyset$, there exists $\varepsilon > 0$ such that $g(\overline{\Omega}) \cap B(y, \varepsilon) = \emptyset$ and so for ε this small,

$$\int_{\Omega} \phi_{\varepsilon}(g(x) - y) \det Dg(x) dx = 0. \blacksquare$$

As noted above, this will end up being $d(g, \Omega, y)$ in this last case where $g^{-1}(y) = \emptyset$.

Lemma 15.1.12 Suppose g, \hat{g} both satisfy Definition 15.1.6. For δ given there, $\delta = \text{dist}(\mathbf{f}(\partial\Omega), \mathbf{y})$,

$$\delta > \|\mathbf{f} - \mathbf{g}\|_{\infty, \bar{\Omega}}, \quad \delta > \|\mathbf{f} - \hat{\mathbf{g}}\|_{\infty, \bar{\Omega}}$$

Then for $t \in [0, 1]$ so does $t\mathbf{g} + (1-t)\hat{\mathbf{g}}$. In particular, $\mathbf{y} \notin (t\mathbf{g} + (1-t)\hat{\mathbf{g}})(\partial\Omega)$. Also $d(\mathbf{f} - \mathbf{y}, \Omega, \mathbf{0}) = d(\mathbf{f}, \Omega, \mathbf{y})$.

Proof: This follows from the fact that $B(\mathbf{y}, \delta)$ in $\|\cdot\|_{\infty, \bar{\Omega}}$ is convex. From the triangle inequality, if $t \in [0, 1]$,

$$\begin{aligned} \|\mathbf{f} - (t\mathbf{g} + (1-t)\hat{\mathbf{g}})\|_{\infty} &\leq t\|\mathbf{f} - \mathbf{g}\|_{\infty} + (1-t)\|\mathbf{f} - \hat{\mathbf{g}}\|_{\infty} \\ &< t\delta + (1-t)\delta = \delta. \end{aligned}$$

If $\|\mathbf{h} - \mathbf{f}\|_{\infty} < \delta$, as was just shown for $\mathbf{h} \equiv t\mathbf{g} + (1-t)\hat{\mathbf{g}}$, then if $x \in \partial\Omega$,

$$\|\mathbf{y} - \mathbf{h}(x)\| \geq \|\mathbf{y} - \mathbf{f}(x)\| - \|\mathbf{h}(x) - \mathbf{f}(x)\| > \text{dist}(\mathbf{f}(\partial\Omega), \mathbf{y}) - \delta \geq \delta - \delta = 0$$

Now consider the last claim. This follows because $\|\mathbf{g} - \mathbf{f}\|_{\infty}$ small is the same as $\|\mathbf{g} - \mathbf{y} - (\mathbf{f} - \mathbf{y})\|_{\infty}$ being small. They are the same. Also, $(\mathbf{g} - \mathbf{y})^{-1}(\mathbf{0}) = \mathbf{g}^{-1}(\mathbf{y})$ and $D\mathbf{g}(x) = D(\mathbf{g} - \mathbf{y})(x)$. ■

Next is an important result on homotopy which is used to show that Definition 15.1.6 is well defined.

Lemma 15.1.13 If \mathbf{h} is in $C^{\infty}(\bar{\Omega} \times [a, b], \mathbb{R}^p)$, and $\mathbf{0} \notin \mathbf{h}(\partial\Omega \times [a, b])$ then for $0 < \varepsilon < \text{dist}(\mathbf{0}, \mathbf{h}(\partial\Omega \times [a, b]))$,

$$t \rightarrow \int_{\Omega} \phi_{\varepsilon}(\mathbf{h}(x, t)) \det D_1 \mathbf{h}(x, t) dx$$

is constant for $t \in [a, b]$. As a special case, $d(\mathbf{f}, \Omega, \mathbf{y})$ is well defined. Also, if $\mathbf{y} \notin \mathbf{f}(\bar{\Omega})$, then $d(\mathbf{f}, \Omega, \mathbf{y}) = 0$.

Proof: By continuity of \mathbf{h} , $\mathbf{h}(\partial\Omega \times [a, b])$ is compact and so is at a positive distance from $\mathbf{0}$. Let $\varepsilon > 0$ be such that for all $t \in [a, b]$,

$$B(\mathbf{0}, \varepsilon) \cap \mathbf{h}(\partial\Omega \times [a, b]) = \emptyset \quad (15.3)$$

Define for $t \in (a, b)$, $H(t) \equiv \int_{\Omega} \phi_{\varepsilon}(\mathbf{h}(x, t)) \det D_1 \mathbf{h}(x, t) dx$. I will show that $H'(t) = 0$ on (a, b) . Then, since H is continuous on $[a, b]$, it will follow from the mean value theorem that $H(t)$ is constant on $[a, b]$. If $t \in (a, b)$,

$$\begin{aligned} H'(t) &= \int_{\Omega} \sum_{\alpha} \phi_{\varepsilon, \alpha}(\mathbf{h}(x, t)) h_{\alpha, t}(x, t) \det D_1 \mathbf{h}(x, t) dx \\ &+ \int_{\Omega} \phi_{\varepsilon}(\mathbf{h}(x, t)) \sum_{\alpha, j} \det D_1(\mathbf{h}(x, t))_{, \alpha j} h_{\alpha, j t} dx \equiv A + B. \end{aligned} \quad (15.4)$$

In this formula, the function \det is considered as a function of the n^2 entries in the $n \times n$ matrix and the $, \alpha j$ represents the derivative with respect to the αj^{th} entry $h_{\alpha, j}$. Now as in the proof of Lemma 7.11.2 on Page 201, $\det D_1(\mathbf{h}(x, t))_{, \alpha j} = (\text{cof } D_1(\mathbf{h}(x, t)))_{\alpha j}$ and so

$$B = \int_{\Omega} \sum_{\alpha} \sum_j \phi_{\varepsilon}(\mathbf{h}(x, t)) (\text{cof } D_1(\mathbf{h}(x, t)))_{\alpha j} h_{\alpha, j t} dx.$$

By hypothesis

$$x \rightarrow \phi_\varepsilon(\mathbf{h}(x, t)) (\text{cof } D_1(\mathbf{h}(x, t)))_{\alpha j} \text{ for } x \in \Omega$$

is in $C_c^\infty(\Omega)$ because if $x \in \partial\Omega$, it follows that for all $t \in [a, b]$, $\mathbf{h}(x, t) \notin B(\mathbf{0}, \varepsilon)$ and so $\phi_\varepsilon(\mathbf{h}(x, t)) = 0$ off some compact set contained in Ω . Therefore, integrate by parts and write

$$\begin{aligned} B &= - \int_\Omega \sum_\alpha \sum_j \frac{\partial}{\partial x_j} (\phi_\varepsilon(\mathbf{h}(x, t))) (\text{cof } D_1(\mathbf{h}(x, t)))_{\alpha j} h_{\alpha, t} dx + \\ &\quad - \int_\Omega \sum_\alpha \sum_j \phi_\varepsilon(\mathbf{h}(x, t)) (\text{cof } D(\mathbf{h}(x, t)))_{\alpha j, j} h_{\alpha, t} dx \end{aligned}$$

The second term equals zero by Lemma 15.1.8. Simplifying the first term yields

$$\begin{aligned} B &= - \int_\Omega \sum_\alpha \sum_j \sum_\beta \phi_{\varepsilon, \beta}(\mathbf{h}(x, t)) h_{\beta, j} h_{\alpha, t} (\text{cof } D_1(\mathbf{h}(x, t)))_{\alpha j} dx \\ &= - \int_\Omega \sum_\alpha \sum_\beta \phi_{\varepsilon, \beta}(\mathbf{h}(x, t)) h_{\alpha, t} \sum_j h_{\beta, j} (\text{cof } D_1(\mathbf{h}(x, t)))_{\alpha j} dx \end{aligned}$$

Now the sum on j is the dot product of the β^{th} row with the α^{th} row of the cofactor matrix which equals zero unless $\beta = \alpha$ because it would be a cofactor expansion of a matrix with two equal rows. When $\beta = \alpha$, the sum on j reduces to $\det(D_1(\mathbf{h}(x, t)))$. Thus B reduces to

$$= - \int_\Omega \sum_\alpha \phi_{\varepsilon, \alpha}(\mathbf{h}(x, t)) h_{\alpha, t} \det(D_1(\mathbf{h}(x, t))) dx$$

Which is the same thing as A , but with the opposite sign. Hence $A + B$ in 15.4 is 0 and $H'(t) = 0$ and so H is a constant on $[a, b]$.

Finally consider the last claim. If $\mathbf{g}, \hat{\mathbf{g}}$ both work in the definition for the degree, then consider $\mathbf{h}(x, t) \equiv t\mathbf{g}(x) + (1-t)\hat{\mathbf{g}}(x) - \mathbf{y}$ for $t \in [0, 1]$. For $x \in \partial\Omega$,

$$\begin{aligned} &|t\mathbf{g}(x) + (1-t)\hat{\mathbf{g}}(x) - \mathbf{y}| \\ &= |t(\mathbf{g}(x) - \mathbf{f}(x)) + (1-t)(\hat{\mathbf{g}}(x) - \mathbf{f}(x)) + \mathbf{f}(x) - \mathbf{y}| \\ &\geq |\mathbf{f}(x) - \mathbf{y}| - |t(\mathbf{g}(x) - \mathbf{f}(x)) + (1-t)(\hat{\mathbf{g}}(x) - \mathbf{f}(x))| \\ &\geq \text{dist}(\mathbf{f}(\partial\Omega), \mathbf{y}) - (t\|\mathbf{g} - \mathbf{f}\|_\infty + (1-t)\|\hat{\mathbf{g}} - \mathbf{f}\|_\infty) \\ &> \text{dist}(\mathbf{f}(\partial\Omega), \mathbf{y}) - (t\delta + (1-t)\delta) = 0 \end{aligned}$$

From Lemma 15.1.12, \mathbf{h} satisfies what is needed for the first part of this lemma. Namely, $\mathbf{0} \notin \mathbf{h}(\partial\Omega \times [0, 1])$. Then from the first part, if $0 < \varepsilon < \text{dist}(\mathbf{0}, \mathbf{h}(\partial\Omega \times [0, 1]))$ and ε is also sufficiently small that the second and last equations hold in what follows,

$$\begin{aligned} d(\mathbf{f}, \Omega, \mathbf{y}) &= \sum \{ \text{sgn}(\det(D\mathbf{g}(x))) : x \in \mathbf{g}^{-1}(\mathbf{y}) \} = \int_\Omega \phi_\varepsilon(\mathbf{h}(x, 1)) \det D_1 \mathbf{h}(x, 1) dx \\ &= \int_\Omega \phi_\varepsilon(\mathbf{h}(x, 0)) \det D_1 \mathbf{h}(x, 0) dx = \sum \{ \text{sgn}(\det(D\hat{\mathbf{g}}(x))) : x \in \hat{\mathbf{g}}^{-1}(\mathbf{y}) \} \blacksquare \end{aligned}$$

15.2 Properties of the Degree

Now that the degree for a continuous function has been defined, it is time to consider properties of the degree. In particular, it is desirable to prove a theorem about homotopy invariance which depends only on continuity considerations.

Theorem 15.2.1 *If \mathbf{h} is in $C(\overline{\Omega} \times [a, b], \mathbb{R}^p)$, and $\mathbf{0} \notin \mathbf{h}(\partial\Omega \times [a, b])$ for each t , then $t \rightarrow d(\mathbf{h}(\cdot, t), \Omega, \mathbf{0})$ is constant for $t \in [a, b]$.*

Proof: Let $0 < \delta = \min |\mathbf{h}(\partial\Omega \times [a, b])|$. By Corollary 15.1.1, there exists $\mathbf{h}_m(\cdot, t) = \sum_{k=0}^m p_k(t) \mathbf{h}(\cdot, t_k)$ for $p_k(t)$ a polynomial in t of degree m such that $p_0(a) = 1$ but $p_k(a) = 0$ if $k \neq 0$ and $p_m(b) = 1$ but $p_k(b) = 0$ if $k \neq m$ and

$$\max_{t \in [a, b]} \|\mathbf{h}_m(\cdot, t) - \mathbf{h}(\cdot, t)\|_{\infty, \overline{\Omega}} < \delta, t_0 = a, t_m = b \quad (15.5)$$

Now replace $\mathbf{h}(\cdot, t_k)$ with $\mathbf{g}_k^m(\cdot) \in C^\infty(\overline{\Omega}, \mathbb{R}^p)$ and $\mathbf{0}$ is a regular value of \mathbf{g}_k^m and let $\mathbf{g}_m(\cdot, t) \equiv \sum_{k=0}^m p_k(t) \mathbf{g}_k^m(\cdot)$ where the functions \mathbf{g}_k^m are close enough to $\mathbf{h}(\cdot, t_k)$ that

$$\max_{t \in [a, b]} \|\mathbf{g}_m(\cdot, t) - \mathbf{h}(\cdot, t)\|_{\infty, \overline{\Omega}} < \delta. \quad (15.6)$$

$\mathbf{g}_m \in C^\infty(\overline{\Omega} \times [a, b]; \mathbb{R}^p)$ because all partial derivatives with respect to either t or \mathbf{x} are continuous. Thus $\mathbf{g}_0^m(\cdot) = \mathbf{g}_m(\cdot, a)$, $\mathbf{g}_m^m(\cdot) = \mathbf{g}_m(\cdot, b)$. Also, from the definition of the degree and Lemma 15.1.13, for small enough ε ,

$$\begin{aligned} d(\mathbf{h}(\cdot, a), \Omega, \mathbf{0}) &= d(\mathbf{g}_0^m(\cdot), \Omega, \mathbf{0}) = \int_{\Omega} \phi_\varepsilon(\mathbf{g}_m(\mathbf{x}, a)) \det D_1 \mathbf{g}_m(\mathbf{x}, a) d\mathbf{x} \\ &= \int_{\Omega} \phi_\varepsilon(\mathbf{g}_m(\mathbf{x}, b)) \det D_1 \mathbf{g}_m(\mathbf{x}, b) d\mathbf{x} = d(\mathbf{g}_m^m(\cdot), \Omega, \mathbf{0}) = d(\mathbf{h}(\cdot, b), \Omega, \mathbf{0}) \end{aligned}$$

Since a, b are arbitrary, this proves the theorem. ■

Now the following theorem is a summary of the main result on properties of the degree.

Theorem 15.2.2 *Definition 15.1.6 is well defined and the degree satisfies the following properties.*

1. (homotopy invariance) *If $\mathbf{h} \in C(\overline{\Omega} \times [0, 1], \mathbb{R}^p)$ and $\mathbf{y}(t) \notin \mathbf{h}(\partial\Omega, t)$ for all $t \in [0, 1]$ where \mathbf{y} is continuous, then*

$$t \rightarrow d(\mathbf{h}(\cdot, t), \Omega, \mathbf{y}(t))$$

is constant for $t \in [0, 1]$.

2. *If $\Omega \supseteq \Omega_1 \cup \Omega_2$ where $\Omega_1 \cap \Omega_2 = \emptyset$, for Ω_i an open set, then if*

$$\mathbf{y} \notin \mathbf{f}(\overline{\Omega} \setminus (\Omega_1 \cup \Omega_2)),$$

then

$$d(\mathbf{f}, \Omega_1, \mathbf{y}) + d(\mathbf{f}, \Omega_2, \mathbf{y}) = d(\mathbf{f}, \Omega, \mathbf{y})$$

3. *$d(I, \Omega, \mathbf{y}) = 1$ if $\mathbf{y} \in \Omega$.*

4. $d(f, \Omega, \cdot)$ is continuous and constant on every connected component of $\mathbb{R}^p \setminus f(\partial\Omega)$.
5. $d(g, \Omega, y) = d(f, \Omega, y)$ if $g|_{\partial\Omega} = f|_{\partial\Omega}$.
6. If $y \notin f(\partial\Omega)$, and if $d(f, \Omega, y) \neq 0$, then there exists $x \in \Omega$ such that $f(x) = y$.

Proof: That the degree is well defined follows from Lemma 15.1.13.

Consider 1., the first property about homotopy. This follows from Theorem 15.2.1 applied to $H(x, t) \equiv h(x, t) - y(t)$.

Consider 2. where $y \notin f(\overline{\Omega} \setminus (\Omega_1 \cup \Omega_2))$. Note that

$$\text{dist}(y, f(\overline{\Omega} \setminus (\Omega_1 \cup \Omega_2))) \leq \text{dist}(y, f(\partial\Omega))$$

Then let g be in $C(\overline{\Omega}; \mathbb{R}^p)$ and

$$\begin{aligned} \|g - f\|_\infty &< \text{dist}(y, f(\overline{\Omega} \setminus (\Omega_1 \cup \Omega_2))) \\ &\leq \min(\text{dist}(y, f(\partial\Omega_1)), \text{dist}(y, f(\partial\Omega_2)), \text{dist}(y, f(\partial\Omega))) \end{aligned}$$

where y is a regular value of g . Then by definition,

$$\begin{aligned} d(f, \Omega, y) &\equiv \sum \{ \det(Dg(x)) : x \in g^{-1}(y) \} \\ &= \sum \{ \det(Dg(x)) : x \in g^{-1}(y), x \in \Omega_1 \} \\ &\quad + \sum \{ \det(Dg(x)) : x \in g^{-1}(y), x \in \Omega_2 \} \\ &\equiv d(f, \Omega_1, y) + d(f, \Omega_2, y) \end{aligned}$$

It is of course obvious that this can be extended by induction to any finite number of disjoint open sets Ω_i .

Note that 3. is obvious because $I(x) = x$ and so if $y \in \Omega$, then $I^{-1}(y) = y$ and $DI(x) = I$ for any x so the definition gives 3.

Now consider 4. Let U be a connected component of $\mathbb{R}^p \setminus f(\partial\Omega)$. This is open as well as connected and arc wise connected by Theorem 3.11.12. Hence, if $u, v \in U$, there is a continuous function $y(t)$ which is in U such that $y(0) = u$ and $y(1) = v$. By homotopy invariance, it follows $d(f, \Omega, y(t))$ is constant. Thus $d(f, \Omega, u) = d(f, \Omega, v)$.

Next consider 5. When $f = g$ on $\partial\Omega$, it follows that if $y \notin f(\partial\Omega)$, then $y \notin f(x) + t(g(x) - f(x))$ for $t \in [0, 1]$ and $x \in \partial\Omega$ so $d(f + t(g - f), \Omega, y)$ is constant for $t \in [0, 1]$ by homotopy invariance in part 1. Therefore, let $t = 0$ and then $t = 1$ to obtain 5.

Claim 6. follows from Lemma 15.1.13 which says that if $y \notin f(\overline{\Omega})$, then $d(f, \Omega, y) = 0$. ■

From the above, there is an easy corollary which gives related properties of the degree.

Corollary 15.2.3 *The following additional properties of the degree are also valid.*

1. If $y \notin f(\overline{\Omega} \setminus \Omega_1)$ and Ω_1 is an open subset of Ω , then $d(f, \Omega, y) = d(f, \Omega_1, y)$.
2. $d(\cdot, \Omega, y)$ is defined and constant on

$$\{g \in C(\overline{\Omega}; \mathbb{R}^p) : \|g - f\|_\infty < r\}$$

where $r = \text{dist}(y, f(\partial\Omega))$.

3. If $y \in f(\Omega)$, $\text{dist}(y, f(\partial\Omega)) \geq \delta$ and $|z - y| < \delta$, then $d(f, \Omega, y) = d(f, \Omega, z)$.

Proof: Consider 1. You can take $\Omega_2 = \emptyset$ in 2 of Theorem 15.2.2 or you can modify the proof of 2 slightly. Consider 2. To verify, let $h(x, t) = f(x) + t(g(x) - f(x))$. Then note that $y \notin h(\partial\Omega, t)$ and use Property 1 of Theorem 15.2.2.

Finally, consider 3. Let $y(t) \equiv (1-t)y + tz$. Then for $x \in \partial\Omega$

$$\begin{aligned} |(1-t)y + tz - f(x)| &= |y - f(x) + t(z - y)| \\ &\geq \delta - t|z - y| > \delta - \delta = 0 \end{aligned}$$

Then by 1 of Theorem 15.2.2, $d(f, \Omega, (1-t)y + tz)$ is constant. When $t = 0$ you get $d(f, \Omega, y)$ and when $t = 1$ you get $d(f, \Omega, z)$. ■

Corollary 15.2.4 Let $h \in C^\infty(\overline{\Omega}, \mathbb{R}^n)$ where Ω is a bounded open set in \mathbb{R}^n and let $y \notin h(\partial\Omega)$. Then $d(h, \Omega, y) = \lim_{\varepsilon \rightarrow 0} \int_{\Omega} \phi_\varepsilon(h(x) - y) \det D h(x) dx$.

Proof: Let $\|\tilde{h} - h\|_{\infty, \Omega} < \delta$ where $0 < \delta < \text{dist}(y, h(\partial\Omega))$ and y is a regular value for \tilde{h} , and $D\tilde{h}(x) = Dh(x)$. Then

$$\begin{aligned} d(h, \Omega, y) &= d(\tilde{h}, \Omega, y) = \lim_{\varepsilon \rightarrow 0} \int_{\Omega} \phi_\varepsilon(\tilde{h}(x) - y) \det D\tilde{h}(x) dx \\ &= \lim_{\varepsilon \rightarrow 0} \int_{\Omega} \phi_\varepsilon(\tilde{h}(x) - y) \det Dh(x) dx \\ &= \lim_{\varepsilon \rightarrow 0} \int_{\Omega} \phi_\varepsilon(h(x) - y) \det Dh(x) dx \end{aligned}$$

because for $h(x, t) = t(h(x) - y) + (1-t)(\tilde{h}(x) - y)$,

$$t \rightarrow \int_{\Omega} \phi_\varepsilon(h(x, t)) \det D_1 h(x, t) dx$$

is constant for $t \in [0, 1]$. ■

15.3 Brouwer Fixed Point Theorem

The degree makes it possible to give a very simple proof of the Brouwer fixed point theorem.

Theorem 15.3.1 (Brouwer fixed point) Let $B = \overline{B(0, r)} \subseteq \mathbb{R}^p$ and let $f : B \rightarrow B$ be continuous. Then there exists a point $x \in B$, such that $f(x) = x$.

Proof: Assume there is no fixed point. Consider $h(x, t) \equiv x - t f(x)$ for $t \in [0, 1]$. Then for $\|x\| = r$, $0 \notin x - t f(x)$, $t \in [0, 1]$. By homotopy invariance, $t \rightarrow d(I - t f, B, 0)$ is constant. But when $t = 0$, this is $d(I, B, 0) = 1 \neq 0$. This is a contradiction so there must be a fixed point after all. ■

You can use standard stuff from Hilbert space to get this the fixed point theorem for a compact convex set. Let K be a closed bounded convex set and let $f : K \rightarrow K$ be continuous. Let P be the projection map onto K as in Problem 10 on Page 152. Then P is

continuous because $|Px - Py| \leq |x - y|$. Recall why this is. From the characterization of the projection map P , $(x - Px, y - Px) \leq 0$ for all $y \in K$. Therefore,

$$(x - Px, Py - Px) \leq 0, (y - Py, Px - Py) \leq 0 \text{ so } (y - Py, Py - Px) \geq 0$$

Hence, subtracting the first from the last,

$$(y - Py - (x - Px), Py - Px) \geq 0$$

consequently,

$$|x - y| |Py - Px| \geq (y - x, Py - Px) \geq |Py - Px|^2$$

and so $|Py - Px| \leq |y - x|$ as claimed.

Now let r be so large that $K \subseteq B(\mathbf{0}, r)$. Then consider $f \circ P$. This map takes $\overline{B(\mathbf{0}, r)} \rightarrow B(\mathbf{0}, r)$. In fact it maps $\overline{B(\mathbf{0}, r)}$ to K . Therefore, being the composition of continuous functions, it is continuous and so has a fixed point in $\overline{B(\mathbf{0}, r)}$ denoted as x . Hence $f(P(x)) = x$. Now, since f maps into K , it follows that $x \in K$. Hence $Px = x$ and so $f(x) = x$. This has proved the following general Brouwer fixed point theorem.

Theorem 15.3.2 *Let $f : K \rightarrow K$ be continuous where K is compact and convex and nonempty, $K \subseteq \mathbb{R}^p$. Then f has a fixed point.*

Definition 15.3.3 *f is a retract of $\overline{B(\mathbf{0}, r)}$ onto $\partial B(\mathbf{0}, r)$ if f is continuous,*

$$f(\overline{B(\mathbf{0}, r)}) \subseteq \partial B(\mathbf{0}, r)$$

and $f(x) = x$ for all $x \in \partial B(\mathbf{0}, r)$.

Theorem 15.3.4 *There does not exist a retract of $\overline{B(\mathbf{0}, r)}$ onto $\partial B(\mathbf{0}, r)$, its boundary.*

Proof: Suppose f were such a retract. Then for all $x \in \partial B(\mathbf{0}, r)$, $f(x) = x$ and so from the properties of the degree, the one which says if two functions agree on $\partial\Omega$, then they have the same degree, $1 = d(I, B(\mathbf{0}, r), \mathbf{0}) = d(f, B(\mathbf{0}, r), \mathbf{0})$ which is clearly impossible because $f^{-1}(\mathbf{0}) = \emptyset$ which implies $d(f, B(\mathbf{0}, r), \mathbf{0}) = 0$. ■

You should now use this theorem to give another proof of the Brouwer fixed point theorem.

15.4 Borsuk's Theorem

In this section is an important theorem which can be used to verify that $d(f, \Omega, y) \neq 0$. This is significant because when this is known, it follows from Theorem 15.2.2 that $f^{-1}(y) \neq \emptyset$. In other words there exists $x \in \Omega$ such that $f(x) = y$.

Definition 15.4.1 *A bounded open set, Ω is symmetric if $-\Omega = \Omega$. A continuous function $f : \overline{\Omega} \rightarrow \mathbb{R}^p$ is odd if $f(-x) = -f(x)$.*

Suppose Ω is symmetric and $g \in C^\infty(\overline{\Omega}; \mathbb{R}^p)$ is an odd map for which $\mathbf{0}$ is a regular value. Then the chain rule implies $Dg(-x) = -Dg(x)$ and so $d(g, \Omega, \mathbf{0})$ must equal an odd integer because if $x \in g^{-1}(\mathbf{0})$, it follows that $-x \in g^{-1}(\mathbf{0})$ also and since $Dg(-x) =$

$Dg(x)$, it follows the overall contribution to the degree from x and $-x$ must be an even integer. Also $0 \in g^{-1}(0)$ and so the degree equals an even integer added to $\text{sgn}(\det Dg(0))$, an odd integer, either -1 or 1 . It seems reasonable to expect that something like this would hold for an arbitrary continuous odd function defined on symmetric Ω . In fact this is the case and this is next. The following lemma is the key result used. This approach is due to Gromes [24]. See also Deimling [12] which is where I found this argument. I think it is one of the cleverest calculus manipulations I have seen.

To get an idea consider the case of $p = 1$. Then Ω is bounded and symmetric and h is odd and in $C^\infty(\overline{\Omega})$. Suppose that $h'(0) \neq 0$. I want to find arbitrarily small ε such that $\hat{h}(x) \equiv h(x) - \varepsilon x^3$ has 0 as a regular value for $x \neq 0$. Let ε be a regular value for $\frac{h(x)}{x^3} \equiv f(x)$ for $x \neq 0$. By Sard's lemma the singular values of $f(x)$ contain no balls so we can take ε as small as desired and have ε a regular value of f . Then at a point where $\hat{h}(x) = 0$, $f(x) = \varepsilon$ and so $\hat{h}(x) + \varepsilon x^3 = x^3 f(x)$. Now differentiate this. $\hat{h}'(x) + 3\varepsilon x^2 = 3x^2 f(x) + x^3 f'(x) = 3x^2 \varepsilon + x^3 f'(x)$ so $\hat{h}'(x) = x^3 f'(x) \neq 0$. This is the motivation for the following process.

The idea is to start with a smooth odd map and approximate it with a smooth odd map which also has 0 a regular value. Note that 0 is a value because $g(0) = -g(0)$.

Process: Suppose $h_0 \in C^\infty(\overline{\Omega}, \mathbb{R}^p)$ is odd and $\det(Dh_0(0)) \neq 0$. Let Ω_k be those points of Ω where $x_k \neq 0$. Here $x \equiv (x_1, \dots, x_p)$. Then $x \rightarrow \frac{h_0(x)}{x_k^3}$ is a smooth map defined on Ω_k so by Sard's lemma, its singular values do not contain $B(0, \eta)$. Therefore, there is y^k with y^k a regular value and $\|y^k\| < \eta$ where $\eta > 0$ is given. Then consider $\hat{h}(x) \equiv h_0(x) - x_k^3 y^k$. I want to argue that 0 is a regular value of \hat{h} on Ω_k . Note that $\frac{h_0(x)}{x_k^3} = y^k$ if and only if $\hat{h}(x) = 0$.

Letting $f(x) \equiv \frac{h_0(x)}{x_k^3} = \frac{\hat{h}(x) + x_k^3 y^k}{x_k^3}$, then $\hat{h}(x) = x_k^3 (f(x) - y^k)$ and $Df(x)$ is invertible at the x of interest, one where $\hat{h}(x) = 0$ and $f(x) - y^k = 0$. Then

$$D\hat{h}(x)(u) = 3x_k^2 \left(\overset{=0}{f(x) - y^k} \right)(u) + x_k^3 Df(x)(u). \quad (15.7)$$

At the point of interest, the first term on the right is 0 and so

$$\det(D\hat{h}(x)) = x_k^3 \det(Df(x)) \neq 0.$$

If 0 is a regular value for h_0 on $\mathcal{U} \subseteq \Omega$, will 0 be a regular value for \hat{h} on \mathcal{U} where \hat{h} is described above? The only points of concern are those $x \in \mathcal{U}$ for which $x_k = 0$ because if $x_k \neq 0$ then $x \in \Omega_k$. But for these points where $x_k = 0$, $\hat{h}(x) = h_0(x)$ and $D\hat{h}(x) = Dh_0(x)$ because $3x_k^2 = 0$ when $x_k = 0$. Thus the new function \hat{h} has 0 a regular value for all $x \in \mathcal{U} \cup \Omega_k$. This **Process** is the basis for the following lemma.

Lemma 15.4.2 Let $h_0 \in C^\infty(\overline{\Omega}, \mathbb{R}^p)$ is odd and $\det(Dh_0(0)) \neq 0$ for Ω a symmetric open set and let $\eta > 0$. Then there are vectors y^k each with $\|y^k\| < \eta$ such that $h(x) \equiv h_0(x) - \sum_{k=1}^p x_k^3 y^k$ has 0 as a regular value.

Proof: Use the above process leading to 15.7 repeatedly. Start with h_0 which has 0 a regular value on $\{0\}$. Then use the process to get $h_1(x) = h_0(x) - y^1 x_1^3$ which has 0 as a regular value on $\{0\} \cup \Omega_1$. Then repeat the process to get $h_2(x) = h_1(x) - y^2 x_2^3$ which has 0 as a regular value on $\{0\} \cup \Omega_1 \cup \Omega_2$. Continue this way and let $h = h_p$ which has 0 a regular value on $\{0\} \cup \Omega_1 \cup \dots \cup \Omega_p = \Omega$. ■

Lemma 15.4.3 *Let $g \in C^\infty(\bar{\Omega}; \mathbb{R}^p)$ be an odd map. Then for every $\varepsilon > 0$, there exists $h \in C^\infty(\bar{\Omega}; \mathbb{R}^p)$ such that h is also an odd map, $\|h - g\|_\infty < \varepsilon$, and 0 is a regular value of h , $0 \notin g(\partial\Omega)$. Here Ω is a symmetric bounded open set. In addition, $d(g, \Omega, 0)$ is an odd integer.*

Proof: In this argument $\eta > 0$ will be a small positive number. Let $h_0(x) = g(x) + \eta x$ where η is sufficiently small but nonzero that $\det Dh_0(0) \neq 0$. See Lemma 8.10.2. Note that h_0 is odd and 0 is a value of h_0 thanks to $h_0(0) = 0$. This has taken care of 0 . However, it is not known whether 0 is a regular value of h_0 because there may be other x where $h_0(x) = 0$. By Lemma 15.4.2, there are vectors y^j with $\|y^k\| \leq \eta$ and 0 is a regular value of $h(x) \equiv h_0(x) - \sum_{j=1}^p y^j x_j^3$. Then

$$\begin{aligned} \|h - g\|_{\infty, \bar{\Omega}} &\leq \max_{x \in \Omega} \left\{ \|\eta x\| + \sum_{k=1}^p \|y^k\| \|x\| \right\} \\ &\leq \eta((p+1) \text{diam}(\Omega)) < \varepsilon < \text{dist}(g(\partial\Omega), 0) \end{aligned}$$

provided η was chosen sufficiently small to begin with.

So what is $d(h, \Omega, 0)$? Since 0 is a regular value and h is odd,

$$h^{-1}(0) = \{x_1, \dots, x_r, -x_1, \dots, -x_r, 0\}.$$

So consider $Dh(x)$ and $Dh(-x)$.

$$\begin{aligned} Dh(-x)u + o(u) &= h(-x+u) - h(-x) = -h(x+(-u)) + h(x) \\ &= -(Dh(x)(-u)) + o(-u) = Dh(x)(u) + o(u) \end{aligned}$$

Hence $Dh(x) = Dh(-x)$ and so the determinants of these two are the same. It follows from the definition that $d(g, \Omega, 0) = d(h, \Omega, 0)$

$$\begin{aligned} &= \sum_{i=1}^r \text{sgn}(\det(Dh(x_i))) + \sum_{i=1}^r \text{sgn}(\det(Dh(-x_i) + \text{sgn}(\det(Dh(0)))))) \\ &= 2m \pm 1 \text{ some integer } m \blacksquare \end{aligned}$$

Theorem 15.4.4 (Borsuk) *Let $f \in C(\bar{\Omega}; \mathbb{R}^p)$ be odd and let Ω be symmetric with $0 \notin f(\partial\Omega)$. Then $d(f, \Omega, 0)$ equals an odd integer.*

Proof: Let ψ_n be a mollifier which is symmetric, $\psi(-x) = \psi(x)$. Also recall that f is the restriction to $\bar{\Omega}$ of a continuous function, still denoted as f which is defined on all of \mathbb{R}^p . Let g be the odd part of this function. That is,

$$g(x) \equiv \frac{1}{2}(f(x) - f(-x)) = f(x) \text{ on } \bar{\Omega}$$

Thus $d(f, \Omega, 0) = d(g, \Omega, 0)$. Then

$$\begin{aligned} g_n(-x) &\equiv g * \psi_n(-x) = \int_{\Omega} g(-x-y) \psi_n(y) dy \\ &= - \int_{\Omega} g(x+y) \psi_n(y) dy = - \int_{\Omega} g(x-(-y)) \psi_n(-y) dy = -g_n(x) \end{aligned}$$

Thus g_n is odd and is infinitely differentiable. Let n be large enough that

$$\|g_n - g\|_{\infty, \bar{\Omega}} < \delta < \text{dist}(f(\partial\Omega), 0) = \text{dist}(g(\partial\Omega), 0)$$

Then by definition of the degree, $d(f, \Omega, 0) = d(g, \Omega, 0) = d(g_n, \Omega, 0)$ and by Lemma 15.4.3 this is an odd integer. \blacksquare

15.5 Some Applications

With Borsuk's theorem it is possible to give relatively easy proofs of some very important and difficult theorems.

Lemma 15.5.1 *Let $g : \overline{B(\mathbf{0}, r)} \subseteq \mathbb{R}^p \rightarrow \mathbb{R}^p$ be one to one and continuous. Then there exists $\delta > 0$ such that $B(g(\mathbf{0}), \delta) \subseteq g(B(\mathbf{0}, r))$.*

Proof: For $t \in [0, 1]$, let $h(x, t) \equiv g(x) - g(-tx)$. Then for $x \in \partial B(\mathbf{0}, r)$, $h(x, t) \neq \mathbf{0}$ because if this were so, the fact g is one to one implies $x = -tx$ and this requires $x = \mathbf{0}$, not the case since $\|x\| = r$. Then $d(h(\cdot, t), B(\mathbf{0}, r), \mathbf{0})$ is constant by Theorem 15.2.1, homotopy invariance. Hence it is an odd integer for all t thanks to Borsuk's theorem, because $h(\cdot, 1)$ is odd. Now let $B(\mathbf{0}, \delta)$ be such that $B(\mathbf{0}, \delta) \cap h(\partial\Omega, 0) = \emptyset$. Then $0 \neq d(h(\cdot, 0), B(\mathbf{0}, r), \mathbf{0}) = d(h(\cdot, 0), B(\mathbf{0}, r), z)$ for $z \in B(\mathbf{0}, \delta)$ because the degree is constant on connected components of $\mathbb{R}^p \setminus h(\partial\Omega, 0)$ by Theorem 15.2.2. Hence $z = h(x, 0) = g(x) - g(\mathbf{0})$ for some $x \in B(\mathbf{0}, r)$. Thus

$$g(B(\mathbf{0}, r)) \supseteq g(\mathbf{0}) + B(\mathbf{0}, \delta) = B(g(\mathbf{0}), \delta). \blacksquare$$

Theorem 15.5.2 (invariance of domain) *Let Ω be any open subset of \mathbb{R}^p and let $f : \Omega \rightarrow \mathbb{R}^p$ be continuous and one to one. Then f maps open subsets of Ω to open sets in \mathbb{R}^p .*

Proof: Let $\overline{B(x_0, r)} \subseteq \Omega$ where f is one to one on $\overline{B(x_0, r)}$. Let g be defined on $B(\mathbf{0}, r)$ given by

$$g(x) \equiv f(x + x_0)$$

Then g satisfies the conditions of Lemma 15.5.1, being one to one and continuous. It follows from that lemma that there exists $\delta > 0$ such that

$$\begin{aligned} f(\Omega) &\supseteq f(B(x_0, r)) = f(x_0 + B(\mathbf{0}, r)) \\ &= g(B(\mathbf{0}, r)) \supseteq g(\mathbf{0}) + B(\mathbf{0}, \delta) \\ &= f(x_0) + B(\mathbf{0}, \delta) = B(f(x_0), \delta) \end{aligned}$$

This shows that for any $x_0 \in \Omega$, $f(x_0)$ is an interior point of $f(\Omega)$ which shows $f(\Omega)$ is open. \blacksquare

Definition 15.5.3 *If $f : U \subseteq \mathbb{R}^p \rightarrow \mathbb{R}^p$ where U is an open set. Then f is locally one to one if for every $x \in U$, there exists $\delta > 0$ such that f is one to one on $B(x, \delta)$.*

Then an examination of the proof of the above theorem shows the following corollary.

Corollary 15.5.4 *In Theorem 15.5.2 it suffices to assume f is locally one to one.*

With the above, one gets easily the following amazing result. It is something which is clear for linear maps but this is a statement about continuous maps.

Corollary 15.5.5 *If $p > m$ there does not exist a continuous one to one map from \mathbb{R}^p to \mathbb{R}^m .*

Proof: Suppose not and let \mathbf{f} be such a continuous map, $\mathbf{f}(\mathbf{x}) \equiv (f_1(\mathbf{x}), \dots, f_m(\mathbf{x}))^T$. Then let $\mathbf{g}(\mathbf{x}) \equiv (f_1(\mathbf{x}), \dots, f_m(\mathbf{x}), 0, \dots, 0)^T$ where there are $p-m$ zeros added in. Then \mathbf{g} is a one to one continuous map from \mathbb{R}^p to \mathbb{R}^p and so $\mathbf{g}(\mathbb{R}^p)$ would have to be open from the invariance of domain theorem and this is not the case. ■

Corollary 15.5.6 *Let $\mathbf{f} : \mathbb{R}^p \rightarrow \mathbb{R}^p$ and $\lim_{|\mathbf{x}| \rightarrow \infty} |\mathbf{f}(\mathbf{x})| = \infty$ where \mathbf{f} is locally one to one and continuous. Then \mathbf{f} maps \mathbb{R}^p onto \mathbb{R}^p .*

Proof: By the invariance of domain theorem, $\mathbf{f}(\mathbb{R}^p)$ is an open set. It is also true that $\mathbf{f}(\mathbb{R}^p)$ is a closed set. Here is why. If $\mathbf{f}(\mathbf{x}_k) \rightarrow \mathbf{y}$, the growth condition ensures that $\{\mathbf{x}_k\}$ is a bounded sequence. Taking a subsequence which converges to $\mathbf{x} \in \mathbb{R}^p$ and using the continuity of \mathbf{f} , it follows $\mathbf{f}(\mathbf{x}) = \mathbf{y}$. Thus $\mathbf{f}(\mathbb{R}^p)$ is both open and closed which implies \mathbf{f} must be an onto map since otherwise, \mathbb{R}^p would not be connected. ■

The proofs of the next two theorems make use of the Tietze extension theorem, Theorem 5.8.5.

Theorem 15.5.7 *Let Ω be a symmetric open set in \mathbb{R}^p such that $\mathbf{0} \in \Omega$ and let $\mathbf{f} : \partial\Omega \rightarrow V$ be continuous where V is an m dimensional subspace of \mathbb{R}^p , $m < p$. Then $\mathbf{f}(-\mathbf{x}) = \mathbf{f}(\mathbf{x})$ for some $\mathbf{x} \in \partial\Omega$.*

Proof: You could reduce to the case where $V = \mathbb{R}^m$ if desired. Suppose not. Using the Tietze extension theorem on components of the function, extend \mathbf{f} to all of \mathbb{R}^p , $\mathbf{f}(\overline{\Omega}) \subseteq V$. (Here the extended function is also denoted by \mathbf{f} .) Let $\mathbf{g}(\mathbf{x}) = \mathbf{f}(\mathbf{x}) - \mathbf{f}(-\mathbf{x})$. Then $\mathbf{0} \notin \mathbf{g}(\partial\Omega)$ and so for some $r > 0$, $B(\mathbf{0}, r) \subseteq \mathbb{R}^p \setminus \mathbf{g}(\partial\Omega)$. For $\mathbf{z} \in B(\mathbf{0}, r)$, $d(\mathbf{g}, \Omega, \mathbf{z}) = d(\mathbf{g}, \Omega, \mathbf{0}) \neq 0$ because $B(\mathbf{0}, r)$ is contained in a component of $\mathbb{R}^p \setminus \mathbf{g}(\partial\Omega)$ and Borsuk's theorem implies that $d(\mathbf{g}, \Omega, \mathbf{0}) \neq 0$ since \mathbf{g} is odd. Hence $V \supseteq \mathbf{g}(\Omega) \supseteq B(\mathbf{0}, r)$ and this is a contradiction because V is m dimensional. ■

This theorem is called the Borsuk Ulam theorem. Note that it implies there exist two points on opposite sides of the surface of the earth which have the same atmospheric pressure and temperature, assuming the earth is symmetric and that pressure and temperature are continuous functions. The next theorem is an amusing result which is like combing hair. It gives the existence of a “cowlick”.

Theorem 15.5.8 *Let p be odd and let Ω be an open bounded set in \mathbb{R}^p with $\mathbf{0} \in \Omega$. Suppose $\mathbf{f} : \partial\Omega \rightarrow \mathbb{R}^p \setminus \{\mathbf{0}\}$ is continuous. Then for some $\mathbf{x} \in \partial\Omega$ and $\lambda \neq 0$, $\mathbf{f}(\mathbf{x}) = \lambda \mathbf{x}$.*

Proof: Using the Tietze extension theorem, extend \mathbf{f} to all of \mathbb{R}^p . Also denote the extended function by \mathbf{f} . Suppose for all $\mathbf{x} \in \partial\Omega$, $\mathbf{f}(\mathbf{x}) \neq \lambda \mathbf{x}$ for all $\lambda \in \mathbb{R}$. Then

$$\mathbf{0} \notin t\mathbf{f}(\mathbf{x}) + (1-t)\mathbf{x}, \quad (\mathbf{x}, t) \in \partial\Omega \times [0, 1].$$

$$\mathbf{0} \notin t\mathbf{f}(\mathbf{x}) - (1-t)\mathbf{x}, \quad (\mathbf{x}, t) \in \partial\Omega \times [0, 1].$$

Thus there exists a homotopy of \mathbf{f} and I and a homotopy of \mathbf{f} and $-I$. Then by the homotopy invariance of degree,

$$d(\mathbf{f}, \Omega, \mathbf{0}) = d(I, \Omega, \mathbf{0}), \quad d(\mathbf{f}, \Omega, \mathbf{0}) = d(-I, \Omega, \mathbf{0}).$$

But this is impossible because $d(I, \Omega, \mathbf{0}) = 1$ but $d(-I, \Omega, \mathbf{0}) = (-1)^n = -1$. ■

15.6 Product Formula, Separation Theorem

This section is on the product formula for the degree which is used to prove the Jordan separation theorem. To begin with is a significant observation which is used without comment below. Recall that the connected components of an open set are open. The formula is all about the composition of continuous functions.

$$\Omega \xrightarrow{f} f(\Omega) \subseteq \mathbb{R}^p \xrightarrow{g} \mathbb{R}^p$$

Lemma 15.6.1 *Let $\{K_i\}_{i=1}^N, N \leq \infty$ be the connected components of $\mathbb{R}^p \setminus C$ where C is a closed set. Then $\partial K_i \subseteq C$.*

Proof: Since K_i is a connected component of an open set, it is itself open. See Theorem 3.11.12. Thus ∂K_i consists of all limit points of K_i which are not in K_i . Let p be such a point. If it is not in C then it must be in some other K_j which is impossible because these are disjoint open sets. Thus if x is a point in U it cannot be a limit point of V for V disjoint from U . ■

Definition 15.6.2 *Let the connected components of $\mathbb{R}^p \setminus f(\partial\Omega)$ be denoted by K_i . From the properties of the degree listed in Theorem 15.2.2, $d(f, \Omega, \cdot)$ is constant on each of these components. Denote by $d(f, \Omega, K_i)$ the constant value on the component K_i .*

The following is the product formula. Note that if K is an unbounded component of $f(\partial\Omega)^C$, then $d(f, \Omega, y) = 0$ for all $y \in K$ by homotopy invariance and the fact that for large enough $\|y\|$, $f^{-1}(y) = \emptyset$ since $f(\overline{\Omega})$ is compact.

Theorem 15.6.3 (product formula) *Let $\{K_i\}_{i=1}^\infty$ be the bounded components of $\mathbb{R}^p \setminus f(\partial\Omega)$ for $f \in C(\overline{\Omega}; \mathbb{R}^p)$, let $g \in C(\mathbb{R}^p, \mathbb{R}^p)$, and suppose that $y \notin g(f(\partial\Omega))$ or in other words, $g^{-1}(y) \cap f(\partial\Omega) = \emptyset$. Then*

$$d(g \circ f, \Omega, y) = \sum_{i=1}^{\infty} d(f, \Omega, K_i) d(g, K_i, y). \quad (15.8)$$

All but finitely many terms in the sum are zero. If there are no bounded components of $f(\partial\Omega)^C$, then $d(g \circ f, \Omega, y) = 0$.

Proof: The compact set $f(\overline{\Omega}) \cap g^{-1}(y)$ is contained in $\mathbb{R}^p \setminus f(\partial\Omega)$ so $f(\overline{\Omega}) \cap g^{-1}(y)$ is covered by finitely many of the components K_i one of which may be the unbounded component. Since these components are disjoint, the other components fail to intersect $f(\overline{\Omega}) \cap g^{-1}(y)$. Thus, if K_i is one of these others, either it fails to intersect $g^{-1}(y)$ or K_i fails to intersect $f(\overline{\Omega})$. Thus either $d(f, \Omega, K_i) = 0$ because K_i fails to intersect $f(\overline{\Omega})$ or $d(g, K_i, y) = 0$ if K_i fails to intersect $g^{-1}(y)$. Thus the sum is always a finite sum. I am using Theorem 15.2.2, the part which says that if $y \notin h(\overline{\Omega})$, then $d(h, \Omega, y) = 0$. Note that by Lemma 15.6.1 $\partial K_i \subseteq f(\partial\Omega)$ so $g(\partial K_i) \subseteq g(f(\partial\Omega))$ and so $y \notin g(\partial K_i)$ because it is assumed that $y \notin g(f(\partial\Omega))$.

Let \tilde{g} be in $C^\infty(\mathbb{R}^p, \mathbb{R}^p)$ and let $\|g - \tilde{g}\|_{\infty, f(\overline{\Omega})} < \text{dist}(y, g(f(\partial\Omega)))$. Thus, for each of the finitely many K_i intersecting $f(\overline{\Omega}) \cap g^{-1}(y)$,

$$\begin{aligned} d(g, K_i, y) &= d(\tilde{g}, K_i, y) \text{ and} \\ d(g \circ f, \Omega, y) &= d(\tilde{g} \circ f, \Omega, y) \end{aligned} \quad (15.9)$$

By Lemma 15.1.5, there exists \tilde{g} such that y is a regular value of \tilde{g} in addition to 15.9 and $\tilde{g}^{-1}(y) \cap f(\partial\Omega) = \emptyset$. Then $\tilde{g}^{-1}(y)$ is contained in the union of the K_i along with the unbounded component(s) and by Lemma 15.1.5 $\tilde{g}^{-1}(y)$ is countable. As discussed there, $\tilde{g}^{-1}(y) \cap K_i$ is finite if K_i is bounded. Let $\tilde{g}^{-1}(y) \cap K_i = \{x_j^i\}_{j=1}^{m_i}$, $m_i \leq \infty$. m_i could only be ∞ on the unbounded component.

Now use Lemma 15.1.5 again to get \tilde{f} in $C^\infty(\bar{\Omega}; \mathbb{R}^p)$ such that each x_j^i is a regular value of \tilde{f} on Ω and also $\|\tilde{f} - f\|_\infty$ is very small, so small that

$$d(\tilde{g} \circ \tilde{f}, \Omega, y) = d(\tilde{g} \circ f, \Omega, y) = d(g \circ f, \Omega, y)$$

and $d(\tilde{f}, \Omega, x_j^i) = d(f, \Omega, x_j^i)$ for each i, j .

Thus, from the above,

$$\begin{aligned} d(g \circ f, \Omega, y) &= d(\tilde{g} \circ \tilde{f}, \Omega, y), \\ d(\tilde{f}, \Omega, x_j^i) &= d(f, \Omega, x_j^i) = d(f, \Omega, K_i) \\ d(\tilde{g}, K_i, y) &= d(g, K_i, y) \end{aligned}$$

Is y a regular value for $\tilde{g} \circ \tilde{f}$ on Ω ? Suppose $z \in \Omega$ and $y = \tilde{g} \circ \tilde{f}(z)$ so $\tilde{f}(z) \in \tilde{g}^{-1}(y)$. Then $\tilde{f}(z) = x_j^i$ for some i, j and $D\tilde{f}(z)^{-1}$ exists. Hence

$$D(\tilde{g} \circ \tilde{f})(z) = D\tilde{g}(x_j^i) D\tilde{f}(z),$$

both linear transformations invertible. Thus y is a regular value of $\tilde{g} \circ \tilde{f}$ on Ω .

What of x_j^i in K_i where K_i is unbounded? As observed, the sum of $\text{sgn}(\det D\tilde{f}(z))$ for $z \in \tilde{f}^{-1}(x_j^i)$ is $d(\tilde{f}, \Omega, x_j^i)$ and is 0 because the degree is constant on K_i which is unbounded.

From the definition of the degree, the left side of 15.8 $d(g \circ f, \Omega, y)$ equals

$$\sum \left\{ \text{sgn}(\det D\tilde{g}(\tilde{f}(z))) \text{sgn}(\det D\tilde{f}(z)) : z \in \tilde{f}^{-1}(\tilde{g}^{-1}(y)) \right\}$$

The $\tilde{g}^{-1}(y)$ are the x_j^i . Thus the above is of the form

$$= \sum_i \sum_j \sum_{z \in \tilde{f}^{-1}(x_j^i)} \text{sgn}(\det(D\tilde{g}(x_j^i))) \text{sgn}(\det(D\tilde{f}(z)))$$

As mentioned, if $x_j^i \in K_i$ an unbounded component, then

$$\sum_{z \in \tilde{f}^{-1}(x_j^i)} \text{sgn}(\det(D\tilde{g}(x_j^i))) \text{sgn}(\det(D\tilde{f}(z))) = 0$$

and so, it suffices to only consider bounded components in what follows and the sum makes sense because there are finitely many x_j^i in bounded K_i . This also shows that if there are

no bounded components of $f(\partial\Omega)^C$, then $d(g \circ f, \Omega, y) = 0$. Thus $d(g \circ f, \Omega, y)$ equals

$$\begin{aligned} &= \sum_i \sum_j \operatorname{sgn}(\det(D\tilde{g}(x_j^i))) \sum_{z \in \tilde{f}^{-1}(x_j^i)} \operatorname{sgn}(\det(D\tilde{f}(z))) \\ &= \sum_i d(\tilde{g}, K_i, y) d(\tilde{f}, \Omega, K_i) \end{aligned}$$

To explain the last step,

$$\sum_{z \in \tilde{f}^{-1}(x_j^i)} \operatorname{sgn}(\det(D\tilde{f}(z))) \equiv d(\tilde{f}, \Omega, x_j^i) = d(\tilde{f}, \Omega, K_i).$$

This proves the product formula because \tilde{g} and \tilde{f} were chosen close enough to f, g respectively that

$$\sum_i d(\tilde{f}, \Omega, K_i) d(\tilde{g}, K_i, y) = \sum_i d(f, \Omega, K_i) d(g, K_i, y) \blacksquare$$

Before the general Jordan separation theorem, I want to first consider the examples of most interest.

Recall that if a function f is continuous and one to one on a compact set K , then f is a homeomorphism of K and $f(K)$. Also recall that if U is a nonempty open set, the boundary of U , denoted as ∂U and meaning those points x with the property that for all $r > 0$ $B(x, r)$ intersects both U and U^C , is $\bar{U} \setminus U$.

Proposition 15.6.4 *Let C be a compact set and let $f : C \rightarrow D \subseteq \mathbb{R}^p$, $p \geq 2$ be one to one and continuous so that C and $f(C) \equiv D$ are homeomorphic. Suppose C^C has only one connected component so C^C is connected. Then D^C also has only one component.*

Proof: Extend f , using the Tietze extension theorem on its entries to all of \mathbb{R}^p and let g be an extension of f^{-1} to all of \mathbb{R}^p . Suppose D^C has a bounded component K . Then from Lemma 15.6.1, $\partial K \subseteq D$, $g(\partial K) \subseteq g(D) = C$. It follows that $d(f \circ g, K, z) = 1$ where $z \in K$ because on ∂K , $f \circ g = id$.

If $z \in K$, then $z \neq f \circ g(k)$ for any $k \in \partial K$ because $f \circ g = id$ on $\partial K \subseteq C$, this by Lemma 15.6.1. Then $g(\partial K)^C \supseteq C^C$. If Q is a bounded component of $g(\partial K)^C$ then if Q contains a point of C^C it follows that C^C is connected, has no points of C and hence no points of $g(\partial K)$ so $Q \supseteq C^C$ and Q is not bounded after all. Thus $g(\partial K)^C$ has no bounded components. Then from the product formula Theorem 15.6.3, $d(f \circ g, K, z) = 0$ which is a contradiction. Thus there is no bounded component of D^C . \blacksquare

This says that if a compact set H fails to separate \mathbb{R}^p and if f is continuous and one to one, then also $f(H)$ fails to separate \mathbb{R}^p .

It is obvious that the unit sphere S^{p-1} divides \mathbb{R}^p into two disjoint open sets, the inside and the outside. The following shows that this also holds for any homeomorphic image of S^{p-1} .

Proposition 15.6.5 *Let B be the ball $B(0, 1)$ with S^{p-1} its boundary, $p \geq 2$. Suppose $f : S^{p-1} \rightarrow C \equiv f(S^{p-1}) \subseteq \mathbb{R}^p$ is a homeomorphism. Then C^C also has exactly two components, one bounded and one unbounded.*

Proof: By Proposition 15.6.4 there is at least one component of $f(\partial B)^C$ called K since it is clear that $(S^{p-1})^C$ is not connected. Let f denote the extension of f to all of \mathbb{R}^p and let $g = f^{-1}$ on $f(\partial B)$ where g is also extended using the Tietze extension theorem to all of \mathbb{R}^p . Let H be the unbounded component of $\mathbb{R}^p \setminus S^{p-1}$.

From Lemma 15.6.1, $\partial K \subseteq f(\partial B)$ so $g(\partial K) \subseteq \partial B$. Also,

$$f \circ g(\partial K) \subseteq f \circ g(f(\partial B)) = f(\partial B).$$

Recall that K has no points in $f(\partial B)$ so if $p \in K$, then p cannot be in $f(\partial B)$ and consequently p cannot be in $f \circ g(\partial K)$ either. Summarizing this,

$$\partial K \subseteq f(\partial B), g(\partial K) \subseteq \partial B, f \circ g(\partial K) \cap K = \emptyset$$

Then picking $p \in K$, by the product rule,

$$1 = d(id, K, p) = d(f \circ g, K, p) = \sum_i d(g, K, Q_i) d(f, Q_i, p)$$

where here the Q_i are the bounded components of $(g(\partial K))^C$. These are maximal open connected sets in \mathbb{R}^p . Recall $g(\partial K) \subseteq \partial B$. If Q_i has a point of H , then H would be connected and contain no points of $g(\partial K)$ and so H would be contained in Q_i which does not happen because Q_i is bounded. Thus $Q_i \subseteq \bar{B}$ but also Q_i is open and so it must be contained in B . Now B is connected and open and contains no points of $g(\partial K)$ because it contains no points of ∂B which is a larger set than $g(\partial K)$ and so in fact $Q_i = B$ and there is only one term in the above sum. Thus, from properties of the degree,

$$\begin{aligned} 1 &= d(id, K, p) = d(f \circ g, K, p) = d(g, K, B) d(f, B, p) \\ &= d(g, K, 0) d(f, B, K) = d(g \circ f, B, 0) \end{aligned}$$

so by the product rule there is no more than one bounded component of $f(\partial B)^C$ the K just mentioned. To emphasize this, if you had bounded components K_i of $f(\partial B)^C$, $i \leq m \leq \infty$. Then $1 = d(g, K_i, 0) d(f, B, K_i) = d(g \circ f, B, 0)$, but then, by the product rule, you would have for $K \equiv K_0$, $1 = d(g \circ f, B, 0) = \sum_{k=0}^m d(g, K_i, 0) d(f, B, K_i) \stackrel{=1}{=} m + 1$. Thus there is exactly one bounded component of $f(\partial B)^C$. ■

A repeat of the above proof yields the following corollary. Replace B with Ω .

Corollary 15.6.6 *Let $\Omega \subseteq \mathbb{R}^p$, $p \geq 2$ be a bounded open connected set such that $\partial\Omega^C$ has two components, a bounded and an unbounded component. Suppose $f : \partial\Omega \rightarrow C \equiv f(\partial\Omega) \subseteq \mathbb{R}^p$ is a homeomorphism. Then C^C also has exactly two components, one bounded and one unbounded.*

As an application, here is a very interesting little result. It has to do with $d(f, \Omega, f(x))$ in the case where f is one to one and Ω is open and connected. You might imagine this should equal 1 or -1 based on one dimensional analogies. Recall a one to one map defined on an interval is either increasing or decreasing. It either preserves or reverses orientation. It is similar in n dimensions and it is a nice application of the Jordan separation theorem and the product formula.

Proposition 15.6.7 *Let Ω be an open connected bounded set in \mathbb{R}^p , $p \geq 2$ such that $\mathbb{R}^p \setminus \partial\Omega$ consists of two connected components. Let $f \in C(\bar{\Omega}; \mathbb{R}^p)$ be continuous and one to one. Then $f(\Omega)$ is the bounded component of $\mathbb{R}^p \setminus f(\partial\Omega)$ and for $y \in f(\Omega)$, $d(f, \Omega, y)$ either equals 1 or -1 .*

Proof: By the Jordan separation theorem, Corollary 15.6.6, $\mathbb{R}^p \setminus f(\partial\Omega)$ consists of two components, a bounded component B and an unbounded component U . Using the Tietze extension theorem, there exists g defined on \mathbb{R}^p such that $g = f^{-1}$ on $f(\bar{\Omega})$. Thus on $\partial\Omega$, $g \circ f = \text{id}$. It follows from this and the product formula that

$$1 = d(\text{id}, \Omega, g(y)) = d(g \circ f, \Omega, g(y)) = d(g, B, g(y)) d(f, \Omega, B)$$

Therefore, $d(f, \Omega, B) \neq 0$ and so for every $z \in B$, it follows $z \in f(\Omega)$. Thus $B \subseteq f(\Omega)$. On the other hand, $f(\Omega)$ cannot have points in both U and B because it is a connected set. Therefore $f(\Omega) \subseteq B$ and this shows $B = f(\Omega)$. Thus $d(f, \Omega, B) = d(f, \Omega, y)$ for each $y \in B$ and the above formula shows this equals either 1 or -1 because the degree is an integer. ■

The one dimensional case also fits into this although it is easier to do by more elementary means. In the case where $n = 1$, the argument is essentially the same. There is one and only one bounded component for $\mathbb{R} \setminus f(\{a, b\})$. This shows how to generalize orientation. It is just the degree. One could use this to describe an orientable manifold without any direct reference to differentiability.

In the case of $f(S^{p-1})$ one wants to verify that this is the boundary of both components, the bounded one and the unbounded one.

Theorem 15.6.8 *Let S^{p-1} be the unit sphere in \mathbb{R}^p , $p \geq 2$. Suppose $\gamma : S^{p-1} \rightarrow \Gamma \subseteq \mathbb{R}^p$ is one to one onto and continuous. Then $\mathbb{R}^p \setminus \Gamma$ consists of two components, a bounded component (called the inside) U_i and an unbounded component (called the outside), U_o . Also the boundary of each of these two components of $\mathbb{R}^p \setminus \Gamma$ is Γ and Γ has empty interior.*

Proof: γ^{-1} is continuous since S^{p-1} is compact and γ is one to one. By the Jordan separation theorem, $\mathbb{R}^p \setminus \Gamma = U_o \cup U_i$ where these on the right are the connected components of the set on the left, both open sets. Only U_i is bounded. Thus $\Gamma \cup U_i \cup U_o = \mathbb{R}^p$. Since both U_i, U_o are open, $\partial U \equiv \bar{U} \setminus U$ for U either U_o or U_i . If $x \in \Gamma$, and is not a limit point of U_i , then there is $B(x, r)$ which contains no points of U_i . Let S be those points x of Γ for which, $B(x, r)$ contains no points of U_i for some $r > 0$. This S is open in Γ . Let $\hat{\Gamma}$ be $\Gamma \setminus S$. Then if $\hat{C} = \gamma^{-1}(\hat{\Gamma})$, it follows that \hat{C} is a closed set in S^{p-1} and is a proper subset of S^{p-1} . It is obvious that taking a relatively open set from S^{p-1} results in a compact set whose complement in \mathbb{R}^p is an open connected set. By Proposition 15.6.4, $\mathbb{R}^p \setminus \hat{\Gamma}$ is also an open connected set. Start with $x \in U_i$ and consider a continuous curve which goes from x to $y \in U_o$ which is contained in $\mathbb{R}^p \setminus \hat{\Gamma}$. Thus the curve contains no points of $\hat{\Gamma}$. However, it must contain points of Γ which can only be in S . The first point of Γ intersected by this curve is a point in \bar{U}_i and so this point of intersection is not in S after all because every ball containing it must contain points of U_i . Thus $S = \emptyset$ and every point of Γ is in \bar{U}_i . Similarly, every point of Γ is in \bar{U}_o . Thus $\Gamma \subseteq \bar{U}_i \setminus U_i$ and $\Gamma \subseteq \bar{U}_o \setminus U_o$. However, if $x \in \bar{U}_i \setminus U_i$, then $x \notin U_o$ because it is a limit point of U_i and so $x \in \Gamma$. It is similar with U_o . Thus $\Gamma = \bar{U}_i \setminus U_i$ and $\Gamma = \bar{U}_o \setminus U_o$. This could not happen if Γ had an interior point. Such a point would be in Γ but would fail to be in either ∂U_i or ∂U_o . ■

When $p = 2$, this theorem is called the Jordan curve theorem.

What if γ maps \bar{B} to \mathbb{R}^p instead of γ only being defined on S^{p-1} ? Obviously, one should be able to say a little more.

Corollary 15.6.9 *Let B be an open ball and let $\gamma : \bar{B} \rightarrow \mathbb{R}^p$ be one to one and continuous. Let U_i, U_o be as in the above theorem, the bounded and unbounded components of $\gamma(\partial B)^C$. Then $U_i = \gamma(B)$.*

Proof: This follows from Proposition 15.6.7.

Note how this yields the invariance of domain theorem. If f is one to one on U an open set, you could consider $\bar{B} \subseteq U$ and then $f(B)$ is the bounded component of $f(\partial B)^C$. You can do this for each ball contained in U . Thus $f(U)$ is open.

15.7 General Jordan Separation Theorem

What follows is the general Jordan separation theorem. First note that if C, D are compact sets and $f : C \rightarrow D$ is a homeomorphism, continuous, one to one and onto, then if C, D are both in \mathbb{R} and if C^C has no bounded components, then C would be a closed interval and so would D . Thus C^C, D^C have the same number of bounded components. In general for \mathbb{R}^p , Proposition 15.6.4 says C^C, D^C both have no bounded components together. The Jordan Separation Theorem shows that C^C, D^C have the same number of bounded components in general.

Lemma 15.7.1 *Let Ω be a bounded open set in \mathbb{R}^p , $f \in C(\bar{\Omega}; \mathbb{R}^p)$, and suppose the sequence $\{\Omega_i\}_{i=1}^\infty$ are disjoint open sets contained in Ω such that*

$$y \notin f(\bar{\Omega} \setminus \cup_{j=1}^\infty \Omega_j)$$

Then $d(f, \Omega, y) = \sum_{j=1}^\infty d(f, \Omega_j, y)$ where the sum has only finitely many terms equal to 0.

Proof: By assumption, the compact set $f^{-1}(y) \equiv \{x \in \bar{\Omega} : f(x) = y\}$ has empty intersection with $\bar{\Omega} \setminus \cup_{j=1}^\infty \Omega_j$ and so this compact set is covered by finitely many of the Ω_j , say $\{\Omega_1, \dots, \Omega_{n-1}\}$ and $y \notin f(\cup_{j=n}^\infty \Omega_j)$. By Theorem 15.2.2 and letting $O = \cup_{j=n}^\infty \Omega_j$,

$$d(f, \Omega, y) = \sum_{j=1}^{n-1} d(f, \Omega_j, y) + d(f, O, y) = \sum_{j=1}^\infty d(f, \Omega_j, y)$$

because $d(f, O, y) = 0$ as is $d(f, \Omega_j, y)$ for every $j \geq n$. ■

Theorem 15.7.2 (Jordan separation theorem) *Let f be a homeomorphism of C and $f(C) \equiv D$ where C is a compact set in \mathbb{R}^p . Then $\mathbb{R}^p \setminus C$ and $\mathbb{R}^p \setminus D$ have the same number of connected components.*

Proof: If either C or D has no bounded components, then so does the other, this from Proposition 15.6.4. Let f denote a Tietze extension of f to all of \mathbb{R}^p and let g be a Tietze extension of f^{-1} to all of \mathbb{R}^p . Let the bounded components of C^C be $\{J_r\}_{r=1}^n \equiv \mathcal{J}$ and let the bounded components of D^C be $\{K_s\}_{s=1}^m \equiv \mathcal{K}$, $n, m \leq \infty$. If both are ∞ then we consider the theorem proved. Assume one of n, m is less than ∞ . Pick $x_r \in J_r$ and $y_s \in K_s$. By Lemma 15.6.1, $\partial K_s \subseteq D$ and so $g(\partial K_s) \subseteq g(D) = C$. $f \circ g(\partial K_s) \subseteq f(C) = D$ and K_s is a component of D^C and so $y_s \notin f \circ g(\partial K_s)$. Then from the definition of the degree and its properties along with the product formula,

$$1 = d(f \circ g, K_s, y_s) = \sum_j d(g, K_s, Q_j) d(f, Q_j, y_s) \quad (15.10)$$

where the Q_j are the bounded components of $g(\partial K_s)^C$. If the unbounded component of C^C is U , then considering Q_j , it can't have any point of U . This is because U has no points

of $g(\partial K_s)$ a smaller set than C and so $Q_j \cup U$ would be connected, open, and contained in $g(\partial K_s)^C$ so it would equal Q_j resulting in Q_j not being bounded after all. Could Q_j intersect some J_r ? If it does, then $J_r \subseteq Q_j$ because J_r is connected and does not intersect $g(\partial K_s)^C$. Consider $f(\bar{Q}_j \setminus \cup \mathcal{J}_j)$ where \mathcal{J}_j are the components J_r contained in Q_j . Is $y_s \in f(\bar{Q}_j \setminus \cup \mathcal{J}_j)$? From Lemma 15.6.1, $\partial Q_j \subseteq g(\partial K_s) \subseteq C$ so $f(\partial Q_j) \subseteq \partial K_s$ and so $y_s \notin f(\partial Q_j)$. Suppose $y_s = f(z)$ where $z \in Q_j$. If z is not in any of the J_r but is in \bar{Q}_j then $z \in C$ so $f(z) = y_s \in D$. But y_s is in K_s a component of D^C so this is impossible. Hence z is in one of the J_r and so this J_r is in \mathcal{J}_j . Therefore, $y_s \notin f(\bar{Q}_j \setminus \cup \mathcal{J}_j)$ and so we can apply Lemma 15.7.1 in 15.10. First note that if $J_r \in \mathcal{J}_j$ then $d(g, K_s, Q_j) = d(g, K_s, J_r)$

$$1 = d(f \circ g, K_s, y_s) = \sum_j d(g, K_s, Q_j) d(f, Q_j, y_s) = \sum_j \sum_{J \in \mathcal{J}_j} d(g, K_s, J) d(f, J, y_s)$$

Since the Q_j cover at least C^C , it follows that each J intersects some Q_j and from the above is contained in Q_j . Thus the \mathcal{J}_j cover $\cup \mathcal{J}$. Therefore, the above equals

$$= \sum_{J \in \mathcal{J}} d(g, K_s, J) d(f, J, y_s) = \sum_{r=1}^n d(g, K_s, J_r) d(f, J_r, K_s)$$

where \mathcal{J} is the set of components of C^C . Recall that in the product formula the sums are finite. Then adding over s , it follows

$$m = \sum_{s=1}^m \sum_{r=1}^n d(g, K_s, J_r) d(f, J_r, K_s)$$

However, we could do the same thing in the other order starting with components in C^C and obtain

$$1 = \sum_{s=1}^m d(g, K_s, J_r) d(f, J_r, K_s)$$

and then summing over r ,

$$n = \sum_{r=1}^n \sum_{s=1}^m d(g, K_s, J_r) d(f, J_r, K_s) = \sum_{s=1}^m \sum_{r=1}^n d(g, K_s, J_r) d(f, J_r, K_s) = m. \blacksquare$$

15.8 Uniqueness of the Degree

I am mainly interested in the topological theorems which can be proved using the above topological degree. To me this justifies its importance. Nevertheless, there are other methods for finding the degree which are based more directly on topological considerations and algebra. These other methods are older than the presentation given here. Nevertheless if the degree satisfies the properties of the degree given in Theorem 15.2.2 along with the following condition, then this is sufficient to determine the degree.

Condition 15.8.1 Let $f : \overline{B(w, R)} \rightarrow \mathbb{R}^p$ be such that $f^{-1}(f(w)) = \{w\}$ and suppose $Df(w)$ is invertible. Then $d(f, B(w, R), f(w)) = \text{sgn}(\det(Df(w)))$.

This follows from a repeat of the arguments which led to the degree in the above. Homotopy invariance and the properties of Theorem 15.2.2 can be used to get the same definition of the degree for continuous functions given in the above. From this all the rest

followed. In an appendix to my book “Linear Algebra and Analysis” such an approach to the degree based on algebra is given and it verifies the above condition. Thus this other approach based on homology gives the same degree function. Also, the above condition will end up following from Theorem 15.2.2 and by insisting that if $s(x) = \hat{x}$ where \hat{x} has two components switched so it corresponds to that elementary matrix then the degree is -1 , this will suffice with the other properties to show the above condition. This process is followed in that other approach to the degree. That something more is required follows because the degree also keeps track of orientation.

15.9 Exercises

1. Show that if y_1, \dots, y_r in $\mathbb{R}^p \setminus f(\partial\Omega)$, then if \tilde{f} has the property that

$$\|\tilde{f} - f\|_\infty < \min_{i \leq r} \text{dist}(y_i, f(\partial\Omega)),$$

then $d(f, \Omega, y_i) = d(\tilde{f}, \Omega, y_i)$ for each y_i . **Hint:** Consider for

$$t \in [0, 1], f(x) + t(\tilde{f}(x) - f(x)) - y_i$$

and homotopy invariance.

2. Show the Brouwer fixed point theorem is equivalent to the nonexistence of a continuous retraction onto the boundary of $B(\mathbf{0}, r)$.
3. Give a version of Proposition 15.6.7 which is valid for the case where $n = 1$.
4. It was shown that if $\lim_{|x| \rightarrow \infty} |f(x)| = \infty$, $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ is locally one to one and continuous, then f maps \mathbb{R}^n onto \mathbb{R}^n . Suppose you have $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$ where f is one to one, continuous, and $\lim_{|x| \rightarrow \infty} |f(x)| = \infty$, $m < n$. Show that f cannot be onto.
5. Can there exist a one to one onto continuous map, f which takes the unit interval to the unit disk?
6. Let $m < n$ and let $B_m(\mathbf{0}, r)$ be the ball in \mathbb{R}^m and $B_n(\mathbf{0}, r)$ be the ball in \mathbb{R}^n . Show that there is no one to one continuous map from $\overline{B_m(\mathbf{0}, r)}$ to $\overline{B_n(\mathbf{0}, r)}$. **Hint:** It is like the above problem.
7. Consider the unit disk, $\{(x, y) : x^2 + y^2 \leq 1\} \equiv D$ and the annulus

$$\left\{ (x, y) : \frac{1}{2} \leq x^2 + y^2 \leq 1 \right\} \equiv A.$$

Is it possible there exists a one to one onto continuous map f such that $f(D) = A$? Thus D has no holes and A is really like D but with one hole punched out. Can you generalize to different numbers of holes? **Hint:** Consider the invariance of domain theorem. The interior of D would need to be mapped to the interior of A . Where do the points of the boundary of A come from? Consider Theorem 3.11.3.

8. Suppose C is a compact set in \mathbb{R}^n which has empty interior and $f : C \rightarrow \Gamma \subseteq \mathbb{R}^n$ is one to one onto and continuous with continuous inverse. Could Γ have nonempty interior? Show also that if f is one to one and onto Γ then if it is continuous, so is f^{-1} .
9. Let K be a nonempty closed and convex subset of \mathbb{R}^p . Recall K is convex means that if $x, y \in K$, then for all $t \in [0, 1]$, $tx + (1-t)y \in K$. Show that if $x \in \mathbb{R}^p$ there exists a unique $z \in K$ such that $|x - z| = \min\{|x - y| : y \in K\}$. This z will be denoted as Px . **Hint:** First note you do not know K is compact. Establish the parallelogram identity if you have not already done so,

$$|u - v|^2 + |u + v|^2 = 2|u|^2 + 2|v|^2.$$

Then let $\{z_k\}$ be a minimizing sequence,

$$\lim_{k \rightarrow \infty} |z_k - x|^2 = \inf\{|x - y|^2 : y \in K\} \equiv \lambda.$$

Using convexity, explain why

$$\left| \frac{z_k - z_m}{2} \right|^2 + \left| x - \frac{z_k + z_m}{2} \right|^2 = 2 \left| \frac{x - z_k}{2} \right|^2 + 2 \left| \frac{x - z_m}{2} \right|^2$$

and then use this to argue $\{z_k\}$ is a Cauchy sequence. Then if z_i works for $i = 1, 2$, consider $(z_1 + z_2)/2$ to get a contradiction.

10. In Problem 9 show that Px satisfies the following variational inequality. $(x - Px) \cdot (y - Px) \leq 0$ for all $y \in K$. Then show that $|Px_1 - Px_2| \leq |x_1 - x_2|$. **Hint:** For the first part note that if $y \in K$, the function

$$t \rightarrow |x - (Px + t(y - Px))|^2$$

achieves its minimum on $[0, 1]$ at $t = 0$. For the second part,

$$(x_1 - Px_1) \cdot (Px_2 - Px_1) \leq 0, (x_2 - Px_2) \cdot (Px_1 - Px_2) \leq 0.$$

Explain why

$$(x_2 - Px_2 - (x_1 - Px_1)) \cdot (Px_2 - Px_1) \geq 0$$

and then use a some manipulations and the Cauchy Schwarz inequality to get the desired inequality.

11. Suppose D is a set which is homeomorphic to $\overline{B(0, 1)}$. This means there exists a continuous one to one map, h such that $h(\overline{B(0, 1)}) = D$ such that h^{-1} is also one to one. Show that if f is a continuous function which maps D to D then f has a fixed point. Now show that it suffices to say that h is one to one and continuous. In this case the continuity of h^{-1} is automatic. Sets which have the property that continuous functions taking the set to itself have at least one fixed point are said to have the fixed point property. Work Problem 7 using this notion of fixed point property. What about a solid ball and a donut? Could these be homeomorphic?
12. Using the definition of the derivative and the Vitali covering theorem, show that if $f \in C^1(\overline{U}, \mathbb{R}^n)$ and ∂U has n dimensional measure zero then $f(\partial U)$ also has measure zero. (This problem has little to do with this chapter. It is a review.)

13. Suppose Ω is any open bounded subset of \mathbb{R}^n which contains $\mathbf{0}$ and that $\mathbf{f} : \overline{\Omega} \rightarrow \mathbb{R}^n$ is continuous with the property that $\mathbf{f}(\mathbf{x}) \cdot \mathbf{x} \geq 0$ for all $\mathbf{x} \in \partial\Omega$. Show that then there exists $\mathbf{x} \in \Omega$ such that $\mathbf{f}(\mathbf{x}) = \mathbf{0}$. Give a similar result in the case where the above inequality is replaced with \leq . **Hint:** You might consider the function $\mathbf{h}(t, \mathbf{x}) \equiv t\mathbf{f}(\mathbf{x}) + (1-t)\mathbf{x}$.
14. Suppose Ω is an open set in \mathbb{R}^n containing $\mathbf{0}$ and suppose that $\mathbf{f} : \overline{\Omega} \rightarrow \mathbb{R}^n$ is continuous and $|\mathbf{f}(\mathbf{x})| \leq |\mathbf{x}|$ for all $\mathbf{x} \in \partial\Omega$. Show \mathbf{f} has a fixed point in $\overline{\Omega}$. **Hint:** Consider $\mathbf{h}(t, \mathbf{x}) \equiv t(\mathbf{x} - \mathbf{f}(\mathbf{x})) + (1-t)\mathbf{x}$ for $t \in [0, 1]$. If $t = 1$ and some $\mathbf{x} \in \partial\Omega$ is sent to $\mathbf{0}$, then you are done. Suppose therefore, that no fixed point exists on $\partial\Omega$. Consider $t < 1$ and use the given inequality.
15. Let Ω be an open bounded subset of \mathbb{R}^n and let $\mathbf{f}, \mathbf{g} : \overline{\Omega} \rightarrow \mathbb{R}^n$ both be continuous, $\mathbf{0} \notin \mathbf{f}(\partial\Omega)$, such that $|\mathbf{f}(\mathbf{x})| - |\mathbf{g}(\mathbf{x})| > 0$ for all $\mathbf{x} \in \partial\Omega$. Show that then $d(\mathbf{f} - \mathbf{g}, \Omega, \mathbf{0}) = d(\mathbf{f}, \Omega, \mathbf{0})$. Show that if there exists $\mathbf{x} \in \mathbf{f}^{-1}(\mathbf{0})$, then there exists $\mathbf{x} \in (\mathbf{f} - \mathbf{g})^{-1}(\mathbf{0})$. **Hint:** Consider $\mathbf{h}(t, \mathbf{x}) \equiv (1-t)\mathbf{f}(\mathbf{x}) + t(\mathbf{f}(\mathbf{x}) - \mathbf{g}(\mathbf{x}))$ and argue $\mathbf{0} \notin \mathbf{h}(t, \partial\Omega)$ for $t \in [0, 1]$.
16. Let $f : \mathbb{C} \rightarrow \mathbb{C}$ where \mathbb{C} is the field of complex numbers. Thus f has a real and imaginary part. Letting $z = x + iy$, $f(z) = u(x, y) + iv(x, y)$. Recall that the norm in \mathbb{C} is given by $|x + iy| = \sqrt{x^2 + y^2}$ and this is the usual norm in \mathbb{R}^2 for the ordered pair (x, y) . Thus complex valued functions defined on \mathbb{C} can be considered as \mathbb{R}^2 valued functions defined on some subset of \mathbb{R}^2 . Such a complex function is said to be analytic if the usual definition holds. That is $f'(z) = \lim_{h \rightarrow 0} \frac{f(z+h) - f(z)}{h}$. In other words,

$$f(z+h) = f(z) + f'(z)h + o(h) \quad (15.11)$$

at a point z where the derivative exists. Let $f(z) = z^n$ where n is a positive integer. Thus $z^n = p(x, y) + iq(x, y)$ for p, q suitable polynomials in x and y . Show this function is analytic. Next show that for an analytic function and u and v the real and imaginary parts, the Cauchy Riemann equations hold, $u_x = v_y$, $u_y = -v_x$. In terms of mappings show 15.11 has the form

$$\begin{aligned} & \begin{pmatrix} u(x+h_1, y+h_2) \\ v(x+h_1, y+h_2) \end{pmatrix} \\ &= \begin{pmatrix} u(x, y) \\ v(x, y) \end{pmatrix} + \begin{pmatrix} u_x(x, y) & u_y(x, y) \\ v_x(x, y) & v_y(x, y) \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \end{pmatrix} + o(\mathbf{h}) \\ &= \begin{pmatrix} u(x, y) \\ v(x, y) \end{pmatrix} + \begin{pmatrix} u_x(x, y) & -v_x(x, y) \\ v_x(x, y) & u_x(x, y) \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \end{pmatrix} + o(\mathbf{h}) \end{aligned}$$

where $\mathbf{h} = (h_1, h_2)^T$ and h is given by $h_1 + ih_2$. Thus the determinant of the above matrix is always nonnegative. Letting B_r denote the ball $B(0, r) = B((0, 0), r)$ show $d(f, B_r, \mathbf{0}) = n$ where $f(z) = z^n$. As a mapping on \mathbb{R}^2 , $\mathbf{f}(x, y) = \begin{pmatrix} u(x, y) \\ v(x, y) \end{pmatrix}$. Thus show $d(\mathbf{f}, B_r, \mathbf{0}) = n$. **Hint:** You might consider $g(z) \equiv \prod_{j=1}^n (z - a_j)$ where the a_j are small real distinct numbers and argue that both this function and f are analytic but that $\mathbf{0}$ is a regular value for \mathbf{g} although it is not so for \mathbf{f} . However, for each a_j small but distinct $d(\mathbf{f}, B_r, \mathbf{0}) = d(\mathbf{g}, B_r, \mathbf{0})$.

17. Using Problem 16, prove the fundamental theorem of algebra as follows. Let $p(z)$ be a nonconstant polynomial of degree n , $p(z) = a_n z^n + a_{n-1} z^{n-1} + \cdots$. Show that for large enough r , $|p(z)| > |p(z) - a_n z^n|$ for all $z \in \partial B(0, r)$. Now from Problem 15 you can conclude $d(p, B_r, 0) = d(f, B_r, 0) = n$ where $f(z) = a_n z^n$.
18. Suppose $f : \mathbb{R}^p \rightarrow \mathbb{R}^p$ satisfies $|f(x) - f(y)| \geq \alpha |x - y|$, $\alpha > 0$. Show that f must map \mathbb{R}^p onto \mathbb{R}^p . **Hint:** First show f is one to one. Then use invariance of domain. Next show, using the inequality, that the points not in $f(\mathbb{R}^p)$ must form an open set because if y is such a point, then there can be no sequence $\{f(x_n)\}$ converging to it. Finally recall that \mathbb{R}^p is connected.
19. Suppose D is a nonempty bounded open set in \mathbb{R}^p and suppose $f : D \rightarrow \partial D$ is continuous with $f(x) = x$ for $x \in \partial D$. Show this cannot happen. **Hint:** Let $y \in D$ and note that id and f agree on ∂D . Therefore, from properties of the degree, $d(f, D, y) = d(\text{id}, D, y)$. Explain why this cannot occur.
20. Assume D is a closed ball in \mathbb{R}^p and suppose $f : D \rightarrow D$ is continuous. Use the above problem to conclude f has a fixed point. **Hint:** If no fixed point, let $g(x)$ be the point on ∂D which results from extending the ray starting at $f(x)$ to x . This would be a continuous map from D to ∂D which does not move any point on ∂D . Draw a picture. This may be the easiest proof of the Brouwer fixed point theorem but note how dependent it is on the properties of the degree.
21. Use Corollary 15.6.9 to prove the invariance of domain theorem that if U is open and $f : U \subseteq \mathbb{R}^p \rightarrow \mathbb{R}^p$ is continuous and one to one, then $f(U)$ is open. This was discussed in the chapter but go through the details.

Chapter 16

Hausdorff Measure

16.1 Lipschitz Functions

Definition 16.1.1 A function $f : U \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^m$ is Lipschitz if there is a constant K such that for all $x, y \in U$, $|f(x) - f(y)| \leq K|x - y|$. We assume $U \neq \emptyset$.

In what follows, dt will be used instead of dm in order to make the notation more familiar.

Lemma 16.1.2 Suppose $f : [a, b] \rightarrow \mathbb{R}$ is Lipschitz continuous. Then f' exists a.e., is in $L^1([a, b])$, and $f(x) = f(a) + \int_a^x f'(t) dt$. In fact, the almost everywhere existence of the derivative holds with only the assumption that f is increasing or of bounded variation on finite intervals. If $f : \mathbb{R} \rightarrow \mathbb{R}$ is Lipschitz, then f' is in $L^1_{loc}(\mathbb{R})$ and the above formula holds.

Proof: Let the Lipschitz constant for f be K . Then let $g(x) \equiv 2Kx - f(x)$ and $h(x) \equiv 2K + f(x)$. Then these are both increasing continuous functions. By Theorem 9.7.4 there are Lebesgue Stieltjes measures μ_f, μ_g satisfying $g(d) - g(c) = \mu_g([c, d]) = \mu_g((c, d))$ with a similar relation for μ_h . Also $\mu_g, \mu_h \ll m_1$ and are Borel measures so by the Radon Nikodym theorem, there exist nonnegative Borel measurable functions α, β such that for all $E \subseteq [a, b]$ Borel, $\mu_g(E) = \int_E \alpha dm$, $\mu_h(E) = \int_E \beta dm$. Let $r(x) \equiv \frac{1}{2}(\beta(x) - \alpha(x))$. It follows that $f(x) = f(a) + \int_a^x r(t) dt$. From the fundamental theorem of calculus, it follows that $r(x) = f'(x)$ a.e. Recall why this is: For $x \in (a, b)$,

$$\left| \frac{f(x+h) - f(x)}{h} \right| \leq 2 \frac{1}{2h} \int_{x-h}^{x+h} |r(t) - f(x)| dt$$

which converges to 0 at Lebesgue points. The last claim follows similarly from the Radon Nikodym theorem and its corollaries. ■

Recall that it was shown earlier that the derivative of an increasing function exists a.e. (Theorem 9.13.4.) This says more.

16.2 Lipschitz Functions and Gateaux Derivatives

Recall the Gateaux derivative is $D_v f(x) \equiv \lim_{h \rightarrow 0} \frac{f(x+hv) - f(x)}{h}$. Each of these is a Borel function because they can be obtained as the limit of a sequence $h_n \rightarrow 0$ of continuous functions.

Corollary 16.2.1 Suppose $f : \mathbb{R}^p \rightarrow \mathbb{R}$ is Lipschitz continuous,

$$|f(x) - f(y)| \leq K|x - y|.$$

Then $f(x+v) - f(x) = \int_0^1 D_v f(x+tv) dt$ where the integrand is the Gateaux derivative and also $|D_v f(x+tv)| \leq K|v|$ a.e. Also $\nabla f(x)$ exists off a set of measure zero.

Proof: $t \rightarrow f(x+tv) - f(x) \equiv g(t)$ is Lipschitz, so by the definition of the Gateaux derivative and Lemma 16.1.2, (See Theorem 7.5.2)

$$\begin{aligned} f(x+v) - f(x) &= \int_0^1 g'(t) dt = \int_0^1 \lim_{h \rightarrow 0} \frac{f(x+tv+h|v|(v/|v|)) - f(x+tv)}{h|v|} |v| dt \\ &= \int_0^1 D_{v/|v|} f(x+tv/|v|) |v| dt = \int_0^1 D_v f(x+tv) dt \end{aligned}$$

Letting $\mathbf{x}_p \equiv (x_1, \dots, x_{p-1}, 0)$ and $\mathbf{v} \equiv \mathbf{e}_p$, it follows that for every \mathbf{x}_p , $\frac{\partial}{\partial x_p} f(\mathbf{x}_p, t)$ exists for a.e. t . Thus $\frac{\partial}{\partial x_p} f$ exists off a set of measure zero. It is similar for the other partial derivatives, and so, taking a union of p exceptional sets of measure zero, it follows that $\nabla f(\mathbf{x})$ exists a.e. ■

16.3 Rademacher's Theorem

It turns out that Lipschitz functions on \mathbb{R}^p can be differentiated a.e. This is called Rademacher's theorem. It also can be shown to follow from the Lebesgue theory of differentiation. We denote $D_{\mathbf{v}}f(\mathbf{x})$ the directional derivative of f in the direction \mathbf{v} . Here \mathbf{v} is a unit vector. In the following lemma, notation is abused slightly. The symbol $f(\mathbf{x}+t\mathbf{v})$ will mean $t \rightarrow f(\mathbf{x}+t\mathbf{v})$ and $\frac{d}{dt}f(\mathbf{x}+t\mathbf{v})$ will refer to the derivative of this function of t . It is a good idea to review Theorem 11.11.5 on integration with polar coordinates because this will be used in what follows. I will also denote as dx the symbol $dm_p(\mathbf{x})$ to save space.

Lemma 16.3.1 *Let $u : \mathbb{R}^p \rightarrow \mathbb{R}$ be Lipschitz with Lipschitz constant K . Let*

$$u_n \equiv u * \phi_n \equiv \int u(\mathbf{x} - \mathbf{y}) dm_p(\mathbf{y})$$

where $\{\phi_n\}$ is a mollifier,

$$\phi_n(\mathbf{y}) \equiv n^p \phi(n\mathbf{y}), \int \phi(\mathbf{y}) dm_p(\mathbf{y}) = 1, \phi(\mathbf{y}) \geq 0, \phi \in C_c^\infty(B(\mathbf{0}, 1))$$

Then

$$\nabla u_n(\mathbf{x}) = \nabla u * \phi_n(\mathbf{x}) \quad (16.1)$$

where ∇u is defined almost everywhere according to Proposition 16.3.4. In fact,

$$\int_a^b \frac{\partial u}{\partial x_i}(\mathbf{x} + t\mathbf{e}_i) dt = u(\mathbf{x} + b\mathbf{e}_i) - u(\mathbf{x} + a\mathbf{e}_i) \quad (16.2)$$

and $\left| \frac{\partial u}{\partial x_i} \right| \leq K$ so $|\nabla u(\mathbf{x})| \leq \sqrt{p}K$ for a.e. \mathbf{x} . Also, $u_n(\mathbf{x}) \rightarrow u(\mathbf{x})$ uniformly on \mathbb{R}^p and for a suitable subsequence, still denoted with n , $\nabla u_n(\mathbf{x}) \rightarrow \nabla u(\mathbf{x})$ for a.e. \mathbf{x} .

Proof: To get the existence of the gradient satisfying the condition given in 16.2, apply Proposition 16.3.4 to each variable. Now

$$\begin{aligned} \frac{u_n(\mathbf{x} + h\mathbf{e}_i) - u_n(\mathbf{x})}{h} &= \int_{\mathbb{R}^p} \left(\frac{u(\mathbf{x} + h\mathbf{e}_i - \mathbf{y}) - u(\mathbf{x} - \mathbf{y})}{h} \right) \phi_n(\mathbf{y}) dm_p(\mathbf{y}) \\ &= \int_{B(\mathbf{0}, \frac{1}{n})} \left(\frac{u(\mathbf{x} + h\mathbf{e}_i - \mathbf{y}) - u(\mathbf{x} - \mathbf{y})}{h} \right) \phi_n(\mathbf{y}) dm_p(\mathbf{y}) \\ &= \int_{B(\mathbf{0}, 1)} \left(\frac{u(\mathbf{x} + h\mathbf{e}_i - \mathbf{y}) - u(\mathbf{x} - \mathbf{y})}{h} \right) \phi_n(\mathbf{y}) dm_p(\mathbf{y}) \end{aligned}$$

Now if $\mathbf{x} - \mathbf{y}$ is off a set of measure zero, the above difference quotient converges to $\frac{\partial u}{\partial x_i}(\mathbf{x} - \mathbf{y})$. You just use Proposition 16.3.4 on the i^{th} variable. If h_k is any sequence converging to 0, you can apply the dominated convergence theorem in the above and obtain

$$\frac{\partial u_n(\mathbf{x})}{\partial x_i} = \int_{B(\mathbf{0}, 1)} \frac{\partial u(\mathbf{x} - \mathbf{y})}{\partial x_i} \phi_n(\mathbf{y}) dm_p(\mathbf{y}) = \frac{\partial u}{\partial x_i} * \phi_n(\mathbf{x})$$

This proves 16.1.

$$\begin{aligned} |u_n(\mathbf{x}) - u(\mathbf{x})| &\leq \int_{\mathbb{R}^p} |u(\mathbf{x} - \mathbf{y}) - u(\mathbf{x})| \phi_n(\mathbf{y}) dm_p(\mathbf{y}) \\ &= \int_{B(\mathbf{0}, \frac{1}{n})} |u(\mathbf{x} - \mathbf{y}) - u(\mathbf{x})| \phi_n(\mathbf{y}) dm_p(\mathbf{y}) \end{aligned}$$

by uniform continuity of u coming from the Lipschitz condition, when n is large enough, this is no larger than $\int_{\mathbb{R}^p} \varepsilon \phi_n(\mathbf{y}) dm_p(\mathbf{y}) = \varepsilon$ and so uniform convergence holds.

Now consider the last claim. From the first part,

$$\begin{aligned} |u_{nx_i}(\mathbf{x}) - u_{x_i}(\mathbf{x})| &= \left| \int_{B(\mathbf{0}, \frac{1}{n})} u_{x_i}(\mathbf{x} - \mathbf{y}) \phi_n(\mathbf{y}) dm_p(\mathbf{y}) - u_{x_i}(\mathbf{x}) \right| \\ &= \left| \int_{B(\mathbf{x}, \frac{1}{n})} u_{x_i}(\mathbf{z}) \phi_n(\mathbf{x} - \mathbf{z}) dm_p(\mathbf{z}) - u_{x_i}(\mathbf{x}) \right| \\ |u_{nx_i}(\mathbf{x}) - u_{x_i}(\mathbf{x})| &\leq \int_{\mathbb{R}^p} |u_{x_i}(\mathbf{x} - \mathbf{y}) - u_{x_i}(\mathbf{x})| \phi_n(\mathbf{y}) dm_p(\mathbf{y}) \\ &= \int_{B(\mathbf{0}, \frac{1}{n})} |u_{x_i}(\mathbf{x} - \mathbf{y}) - u_{x_i}(\mathbf{x})| \phi_n(\mathbf{y}) dm_p(\mathbf{y}) \\ &= \int_{B(\mathbf{x}, \frac{1}{n})} |u_{x_i}(\mathbf{z}) - u_{x_i}(\mathbf{x})| \phi_n(\mathbf{x} - \mathbf{z}) dm_p(\mathbf{z}) \\ &\leq n^p \int_{B(\mathbf{x}, \frac{1}{n})} |u_{x_i}(\mathbf{z}) - u_{x_i}(\mathbf{x})| \phi(n(\mathbf{x} - \mathbf{z})) dm_p(\mathbf{z}) \\ &\leq \frac{C}{m_p(B(\mathbf{0}, \frac{1}{n}))} \int_{B(\mathbf{x}, \frac{1}{n})} |u_{x_i}(\mathbf{z}) - u_{x_i}(\mathbf{x})| dm_p(\mathbf{z}) \end{aligned}$$

which converges to 0 for a.e. \mathbf{x} , in fact at any Lebesgue point. This is because u_{x_i} is bounded by K and so is in L^1_{loc} . ■

Note that the above holds just as well if u has values in some \mathbb{R}^m and the same proof would work, replacing $|\cdot|$ with $\|\cdot\|$ or the Euclidean norm $|\cdot|$.

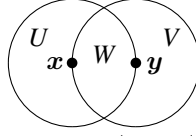
The following lemma gives an interesting inequality due to Morrey. To simplify notation dz will mean $dm_p(\mathbf{z})$.

Lemma 16.3.2 *Let u be a C^1 function on \mathbb{R}^p . Then there exists a constant C , depending only on p such that for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$,*

$$|u(\mathbf{x}) - u(\mathbf{y})| \leq C \left(\int_{B(\mathbf{x}, 2|\mathbf{x} - \mathbf{y}|)} |\nabla u(\mathbf{z})|^q dz \right)^{1/q} \left(|\mathbf{x} - \mathbf{y}|^{(1-p/q)} \right). \quad (16.3)$$

Here $q > p$ and C is some constant depending on p, q .

Proof: In the argument C will be a generic constant which depends on p, q . Consider the following picture.



This is a picture of two balls of radius $r = |x - y|$ in \mathbb{R}^p , U and V having centers at x and y respectively, which intersect in the set W . The center of U is on the boundary of V and the center of V is on the boundary of U as shown in the picture. There exists a constant C , independent of r depending only on p such that $\frac{m_p(W)}{m_p(U)} = \frac{m_p(W)}{m_p(V)} = \frac{1}{C}$. You could compute this constant if you desired but it is not important here.

Then

$$\begin{aligned}
 |u(x) - u(y)| &= \frac{1}{m_p(W)} \int_W |u(x) - u(y)| dz \\
 &\leq \frac{1}{m_p(W)} \int_W |u(x) - u(z)| dz + \frac{1}{m_p(W)} \int_W |u(z) - u(y)| dz \\
 &= \frac{C}{m_p(U)} \left[\int_W |u(x) - u(z)| dz + \int_W |u(z) - u(y)| dz \right] \\
 &\leq \frac{C}{m_p(U)} \left[\int_U |u(x) - u(z)| dz + \int_V |u(y) - u(z)| dz \right] \quad (16.4)
 \end{aligned}$$

Now consider these two terms. Let $q > p$. Consider the first term.

Letting U_0 denote the ball of the same radius as U but with center at 0 .

$$\begin{aligned}
 \frac{1}{m_p(U)} \int_U |u(x) - u(z)| dz &= \frac{1}{m_p(U_0)} \int_{U_0} |u(x) - u(z+x)| dz \\
 &= \frac{1}{m_p(U_0)} \int_{U_0} \left| \int_0^1 \nabla u(x+tz) \cdot z dt \right| dz \leq \frac{1}{m_p(U_0)} \int_0^1 \int_{U_0} |\nabla u(x+tz)| |z| dz dt \\
 &\leq \frac{1}{m_p(U_0)} \int_0^1 \left(\int_{U_0} |\nabla u(x+tz)|^q dz \right)^{1/q} \left(\int_{U_0} |z|^{q/(q-1)} dz \right)^{(q-1)/q} \\
 &= \frac{1}{m_p(U_0)} \int_0^1 \left(\int_{U_0} |\nabla u(x+tz)|^q dz \right)^{1/q} \left(\int_{S^{p-1}} \int_0^r \rho^{q/(q-1)} \rho^{p-1} d\rho d\sigma \right)^{(q-1)/q} \\
 &= C_{pq} \frac{r^{\frac{1}{q}(q-p+pq)}}{m_p(U_0)} \int_0^1 \left(\int_{U_0} |\nabla u(x+tz)|^q dz \right)^{1/q} dt \\
 &= \frac{C_{pq}}{\alpha(p)} \frac{r^{1-\frac{p}{q}}}{m_p(U_0)} \int_0^1 \left(\int_{U_0} |\nabla u(x+tz)|^q dz \right)^{1/q} dt
 \end{aligned}$$

where $C_{pq} = \sigma_{p-1} (S^{p-1})^{(q-1)/q} \left(\frac{q-1}{q-p+pq} \right)^{(q-1)/q}$.

Now estimate the last term.

$$\begin{aligned}
 \int_0^1 \left(\int_{U_0} |\nabla u(x+tz)|^q dz \right)^{1/q} dt &= \int_0^1 \left(\frac{1}{t^p} \int_{tU_0} |\nabla u(x+v)|^q dv \right)^{1/q} dt \\
 &\leq \int_0^1 \frac{1}{t^{p/q}} \left(\int_{U_0} |\nabla u(x+v)|^q dv \right)^{1/q} dt
 \end{aligned}$$

$$= \frac{q}{q-p} \left(\int_{U_0} |\nabla u(x+v)|^q dv \right)^{1/q} = \frac{q}{q-p} \left(\int_U |\nabla u(z)|^q dz \right)^{1/q}$$

Since $q > p$. Thus

$$\frac{1}{m_p(U)} \int_U |u(x) - u(z)| dz \leq C \left(\int_U |\nabla u(z)|^q dz \right)^{1/q} \leq C \left(\int_{B(x, 2|x-y|)} |\nabla u(z)|^q dz \right)^{1/q}$$

and similarly

$$\frac{1}{m_p(V)} \int_V |u(x) - u(z)| dz \leq C \left(\int_{B(x, 2|x-y|)} |\nabla u(z)|^q dz \right)^{1/q}$$

$$\text{From 16.4, } |u(x) - u(y)| \leq C \left(\int_{B(x, 2|x-y|)} |\nabla u(z)|^q dz \right)^{1/q} |x - y|^{1-\frac{p}{q}} \blacksquare$$

Corollary 16.3.3 *Let u be Lipschitz on \mathbb{R}^p with constant K . Then there is a constant C depending only on p, q such that*

$$|u(x) - u(y)| \leq C \left(\int_{B(x, 2|x-y|)} |\nabla u(z)|^q dz \right)^{1/q} (|x - y|^{(1-p/q)}). \quad (16.5)$$

Here $q > p$.

Proof: Let $u_n = u * \phi_n$ where $\{\phi_n\}$ is a mollifier as in Lemma 16.3.1. Then from Lemma 16.3.2, there is a constant depending only on p such that

$$|u_n(x) - u_n(y)| \leq C \left(\int_{B(x, 2|x-y|)} |\nabla u_n(z)|^q dz \right)^{1/q} (|x - y|^{(1-p/q)}).$$

Now $|\nabla u_n| = |\nabla u * \phi_n|$ by Lemma 16.3.1 and this last is bounded. Also, by this lemma, $\nabla u_n(z) \rightarrow \nabla u(z)$ a.e. and $u_n(x) \rightarrow u(x)$ for all x . Therefore, by the dominated convergence theorem, pass to the limit as $n \rightarrow \infty$ and obtain 16.5. \blacksquare

Note you can write 16.5 in the form

$$\begin{aligned} |u(x) - u(y)| &\leq C \left(\frac{1}{|x - y|^p} \int_{B(x, 2|x-y|)} |\nabla u(z)|^q dz \right)^{1/q} |x - y| \\ &= \hat{C} \left(\frac{1}{m_p(B(x, 2|x-y|))} \int_{B(x, 2|x-y|)} |\nabla u(z)|^q dz \right)^{1/q} |x - y| \end{aligned}$$

Before leaving this remarkable formula, note that if you are in any situation where the above formula holds and ∇u exists in some sense and is in L^q , $q > p$, then u would need to be continuous. This is the basis for the Sobolev embedding theorem.

Here is Rademacher's theorem.

Theorem 16.3.4 *Suppose u is Lipschitz with constant K then if x is a point where $\nabla u(x)$ exists,*

$$\begin{aligned} &|u(y) - u(x) - \nabla u(x) \cdot (y - x)| \\ &\leq C \left(\frac{1}{m(B(x, 2|x-y|))} \int_{B(x, 2|x-y|)} |\nabla u(z) - \nabla u(x)|^q dz \right)^{1/q} |x - y|. \end{aligned} \quad (16.6)$$

Also u is differentiable at a.e. x and also

$$u(x + tv) - u(x) = \int_0^t D_v u(x + sv) ds \quad (16.7)$$

Proof: This follows easily from letting $g(\mathbf{y}) \equiv u(\mathbf{y}) - u(\mathbf{x}) - \nabla u(\mathbf{x}) \cdot (\mathbf{y} - \mathbf{x})$. As explained above, $|\nabla u(\mathbf{x})| \leq \sqrt{p}K$ at every point where ∇u exists, the exceptional points being in a set of measure zero. Then $g(\mathbf{x}) = 0$, and $\nabla g(\mathbf{y}) = \nabla u(\mathbf{y}) - \nabla u(\mathbf{x})$ at the points \mathbf{y} where the gradient of g exists. From Corollary 16.3.3,

$$\begin{aligned} & |u(\mathbf{y}) - u(\mathbf{x}) - \nabla u(\mathbf{x}) \cdot (\mathbf{y} - \mathbf{x})| = |g(\mathbf{y})| = |g(\mathbf{y}) - g(\mathbf{x})| \\ & \leq C \left(\int_{B(\mathbf{x}, 2|\mathbf{x} - \mathbf{y}|)} |\nabla u(\mathbf{z}) - \nabla u(\mathbf{x})|^q d\mathbf{z} \right)^{1/q} |\mathbf{x} - \mathbf{y}|^{1 - \frac{p}{q}} \\ & = C \left(\int_{B(\mathbf{x}, 2|\mathbf{x} - \mathbf{y}|)} |\nabla u(\mathbf{z}) - \nabla u(\mathbf{x})|^q d\mathbf{z} \right)^{1/q} \frac{1}{|\mathbf{x} - \mathbf{y}|^{\frac{1}{p}}} |\mathbf{x} - \mathbf{y}| \\ & = C \left(\frac{1}{m(B(\mathbf{x}, 2|\mathbf{x} - \mathbf{y}|))} \int_{B(\mathbf{x}, 2|\mathbf{x} - \mathbf{y}|)} |\nabla u(\mathbf{z}) - \nabla u(\mathbf{x})|^q d\mathbf{z} \right)^{1/q} |\mathbf{x} - \mathbf{y}|. \end{aligned}$$

Now this is no larger than

$$\leq C \left(\frac{1}{m(B(\mathbf{x}, 2|\mathbf{x} - \mathbf{y}|))} \int_{B(\mathbf{x}, 2|\mathbf{x} - \mathbf{y}|)} |\nabla u(\mathbf{z}) - \nabla u(\mathbf{x})| (2\sqrt{p}K)^{q-1} d\mathbf{z} \right)^{1/q} |\mathbf{x} - \mathbf{y}|$$

It follows that at Lebesgue points of ∇u , the above expression is $o(|\mathbf{x} - \mathbf{y}|)$ and so at all such points u is differentiable. As to 16.7, this follows from an application of Lemma 16.1.2 to $f(t) = u(\mathbf{x} + t\mathbf{v})$. ■

Note that for a.e. \mathbf{x} , $D_{\mathbf{v}}u(\mathbf{x}) = \nabla u(\mathbf{x}) \cdot \mathbf{v}$. If you have a line with direction vector \mathbf{v} , does it follow that $Du(\mathbf{x} + t\mathbf{v})$ exists for a.e. t ? We know the directional derivative exists a.e. t but it might not be clear that it is $\nabla u(\mathbf{x}) \cdot \mathbf{v}$.

For $|\mathbf{w}| = 1$, denote the measure of Section 11.11 defined on the unit sphere S^{p-1} as σ . Let $N_{\mathbf{w}}$ be defined as those $t \in [0, \infty)$ for which $D_{\mathbf{w}}u(\mathbf{x} + t\mathbf{w}) \neq \nabla u(\mathbf{x} + t\mathbf{w}) \cdot \mathbf{w}$.

$$B \equiv \{\mathbf{w} \in S^{p-1} : N_{\mathbf{w}} \text{ has positive measure}\}$$

This is contained in the set of points of \mathbb{R}^p where the derivative of $v(\cdot) \equiv u(\mathbf{x} + \cdot)$ fails to exist. Thus from Section 11.11 the measure of this set is $\int_B \int_{N_{\mathbf{w}}} \rho^{n-1} d\rho d\sigma(\mathbf{w})$. This must equal zero from what was just shown about the derivative of the Lipschitz function v existing a.e. and so $\sigma(B) = 0$. The claimed formula follows from this. Thus we obtain the following corollary.

Corollary 16.3.5 *Let u be Lipschitz. Then for any \mathbf{x} and $\mathbf{v} \in S^{p-1} \setminus B_{\mathbf{x}}$ where $\sigma(B_{\mathbf{x}}) = 0$, it follows that for all t ,*

$$u(\mathbf{x} + t\mathbf{v}) - u(\mathbf{x}) = \int_0^t D_{\mathbf{v}}u(\mathbf{x} + s\mathbf{v}) ds = \int_0^t \nabla u(\mathbf{x} + s\mathbf{v}) \cdot \mathbf{v} ds$$

In all of the above, the function u is defined on all of \mathbb{R}^p . However, it is always the case that Lipschitz functions can be extended off a given set. Thus if a Lipschitz function is defined on some set Ω , then it can always be considered the restriction to Ω of a Lipschitz map defined on all of \mathbb{R}^p .

Theorem 16.3.6 *If $h : \Omega \rightarrow \mathbb{R}^m$ is Lipschitz, then there exists $\bar{h} : \mathbb{R}^p \rightarrow \mathbb{R}^m$ which extends h and is also Lipschitz.*

Proof: It suffices to assume $m = 1$ because if this is shown, it may be applied to the components of h to get the desired result. Suppose

$$|h(x) - h(y)| \leq K|x - y|. \quad (16.8)$$

Define

$$\bar{h}(x) \equiv \inf\{h(w) + K|x - w| : w \in \Omega\}. \quad (16.9)$$

If $x \in \Omega$, then for all $w \in \Omega$, $h(w) + K|x - w| \geq h(x)$ by 16.8. This shows $h(x) \leq \bar{h}(x)$. But also you could take $w = x$ in 16.9 which yields $\bar{h}(x) \leq h(x)$. Therefore $\bar{h}(x) = h(x)$ if $x \in \Omega$.

Now suppose $x, y \in \mathbb{R}^p$ and consider $|\bar{h}(x) - \bar{h}(y)|$. Without loss of generality assume $\bar{h}(x) \geq \bar{h}(y)$. (If not, repeat the following argument with x and y interchanged.) Pick $w \in \Omega$ such that $h(w) + K|y - w| - \varepsilon < \bar{h}(y)$. Then

$$\begin{aligned} |\bar{h}(x) - \bar{h}(y)| &= \bar{h}(x) - \bar{h}(y) \leq h(w) + K|x - w| - \\ &\quad [h(w) + K|y - w| - \varepsilon] \leq K|x - y| + \varepsilon. \end{aligned}$$

Since ε is arbitrary, $|\bar{h}(x) - \bar{h}(y)| \leq K|x - y|$ ■

16.4 Weak Derivatives

A related concept is that of weak derivatives. Let $\Omega \subseteq \mathbb{R}^p$ be an open set. A distribution on Ω is defined to be a linear functional on $C_c^\infty(\Omega)$, called the space of test functions. The space of all such linear functionals will be denoted by $\mathcal{D}'(\Omega)$. Actually, more is sometimes done here. One imposes a topology on $C_c^\infty(\Omega)$ making it into a topological vector space, and when this has been done, $\mathcal{D}'(\Omega)$ is defined as the continuous linear maps. To see this, consult the book by Yosida [60] or the book by Rudin [51]. I am ignoring this topology because in practice, one is usually more interested in some other topology which is much less exotic. Thus $\mathcal{D}'(\Omega)$ is an algebraic dual which has nothing to do with topology.

The following is a basic lemma which will be used in what follows. First recall the following definition.

Definition 16.4.1 For Ω an open set in \mathbb{R}^n , $C_c^\infty(\Omega)$ denotes those functions ϕ which are infinitely differentiable and have compact support in Ω . This is a nonempty set of functions by Lemma 12.5.3.

With this definition, the fundamental lemma is as follows.

Lemma 16.4.2 Suppose $f \in L_{loc}^1(\mathbb{R}^n)$ and suppose $\int f\phi dx = 0$ for all $\phi \in C_c^\infty(\mathbb{R}^n)$. Then $f(x) = 0$ a.e. x .

Proof: Without loss of generality f is real-valued. Let $E \equiv \{x : f(x) > \varepsilon\}$ and let $E_m \equiv E \cap B(0, m)$. We show that $m(E_m) = 0$. If not, there exists an open set V , and a compact set K satisfying

$$K \subseteq E_m \subseteq V \subseteq B(0, m), \quad m_p(V \setminus K) < 4^{-1}m(E_m), \quad \int_{V \setminus K} |f| dx < \varepsilon 4^{-1}m_p(E_m).$$

Let H and W be open sets satisfying $K \subseteq H \subseteq \bar{H} \subseteq W \subseteq \bar{W} \subseteq V$ and let $\bar{H} \prec g \prec W$ where the symbol, \prec , has the same meaning as it does in Definition 3.12.3. That is, g equals 1 on

\bar{H} and has compact support contained in W . Then let ϕ_δ be a mollifier and let $h \equiv g * \phi_\delta$ for δ small enough that $K \prec h \prec V$. Thus

$$\begin{aligned} 0 &= \int f h dx = \int_K f dx + \int_{V \setminus K} f h dx \geq \varepsilon m_p(K) - \varepsilon 4^{-1} m_p(E_m) \\ &\geq \varepsilon (m_p(E_m) - 4^{-1} m_p(E_m)) - \varepsilon 4^{-1} m_p(E_m) \geq 2^{-1} \varepsilon m_p(E_m). \end{aligned}$$

Therefore, $m_p(E_m) = 0$, a contradiction. Thus $m_p(E) \leq \sum_{m=1}^{\infty} m_p(E_m) = 0$ and so, since $\varepsilon > 0$ is arbitrary, $m_p(\{x : f(x) > 0\}) = 0$. Similarly $m(\{x : f(x) < 0\}) = 0$. If f is complex valued, the above applies to the real and imaginary parts. ■

Example: The space $L^1_{loc}(\Omega)$ may be considered as a subset of $\mathcal{D}'(\Omega)$ as follows. $f(\phi) \equiv \int_\Omega f(x) \phi(x) dx$ for all $\phi \in C_c^\infty(\Omega)$. Recall that $f \in L^1_{loc}(\Omega)$ if $f \chi_K \in L^1(\Omega)$ whenever K is compact.

This is well defined thanks to Lemma 16.4.2.

Example: $\delta_x \in \mathcal{D}'(\Omega)$ where $\delta_x(\phi) \equiv \phi(x)$.

It will be observed from the above two examples and a little thought that $\mathcal{D}'(\Omega)$ is truly enormous. We shall define the derivative of a distribution in such a way that it agrees with the usual notion of a derivative on those distributions which are also continuously differentiable functions. With this in mind, let f be the restriction to the open set Ω of a smooth function defined on \mathbb{R}^p . Then $D_{x_i} f$ makes sense and for $\phi \in C_c^\infty(\Omega)$

$$D_{x_i} f(\phi) \equiv \int_\Omega D_{x_i} f(x) \phi(x) dx = - \int_\Omega f D_{x_i} \phi dx = -f(D_{x_i} \phi).$$

Motivated by this, here is the definition of a weak derivative.

Definition 16.4.3 For $T \in \mathcal{D}'(\Omega)$, $D_{x_i} T(\phi) \equiv -T(D_{x_i} \phi)$.

One can continue taking derivatives indefinitely. Thus, $D_{x_i x_j} T \equiv D_{x_i}(D_{x_j} T)$ and it is clear that all mixed partial derivatives are equal because this holds for the functions in $C_c^\infty(\Omega)$. Thus one can differentiate virtually anything, even functions that may be discontinuous everywhere. However the notion of “derivative” is very weak, hence the name, “weak derivatives”.

Example: Let $\Omega = \mathbb{R}$ and let $H(x) \equiv \begin{cases} 1 & \text{if } x \geq 0, \\ 0 & \text{if } x < 0. \end{cases}$ Then

$$DH(\phi) = - \int H(x) \phi'(x) dx = \phi(0) = \delta_0(\phi).$$

Note that in this example, DH is not a function.

What happens when Df is a function?

Theorem 16.4.4 Let $\Omega = (a, b)$ and suppose that f and Df are both in $L^1(a, b)$. Then f is equal to a continuous function a.e., still denoted by f and $f(x) = f(a) + \int_a^x Df(t) dt$.

Proof: Consider $f - \int_a^{(\cdot)} Df(t) dt \equiv T$. Is this function equal to some constant a.e.? Let $\phi \in C_c^\infty(a, b)$. By Fubini's theorem, $DT(\phi) \equiv$

$$\int_a^b \left(f(x) - \int_a^x Df(t) dt \right) \phi'(x) dx = \int_a^b f(x) \phi'(x) dx - \int_a^b \int_a^x Df(t) \phi'(x) dt dx$$

$$\begin{aligned}
&\equiv -Df(\phi) - \int_a^b Df(t) \int_t^b \phi'(x) dx dt = -Df(\phi) + \int_a^b Df(t) \phi(t) dt \\
&= -Df(\phi) + Df(\phi) = 0
\end{aligned}$$

Thus the theorem is proved if it is shown that whenever $DT = 0$, it follows that T is a constant. This is the following lemma.

Lemma 16.4.5 *Let $T \in \mathcal{D}^*(a, b)$ and suppose $DT = 0$. Then there exists a constant C such that $T(\phi) = \int_a^b C \phi dx$.*

Proof: $T(D\phi) = 0$ for all $\phi \in C_c^\infty(a, b)$ from the definition of $DT = 0$. Let $\phi_0 \in C_c^\infty(a, b)$, $\int_a^b \phi_0(x) dx = 1$, and let

$$\psi_\phi(x) = \int_a^x [\phi(t) - \left(\int_a^b \phi(y) dy \right) \phi_0(t)] dt$$

for $\phi \in C_c^\infty(a, b)$. Thus $\psi_\phi \in C_c^\infty(a, b)$ and $D\psi_\phi = \phi - \left(\int_a^b \phi(y) dy \right) \phi_0$. Therefore, $\phi = D\psi_\phi + \left(\int_a^b \phi(y) dy \right) \phi_0$, so $T(\phi) = T(D\psi_\phi) + \left(\int_a^b \phi(y) dy \right) T(\phi_0) = \int_a^b T(\phi_0) \phi(y) dy$. Let $C = T\phi_0$. ■

It follows from this lemma that $f(x) - \int_a^x Df(t) dt = C$ for some constant which we denote as $f(a)$ so that $f(x) = f(a) + \int_a^x Df(t) dt$. ■

Theorem 16.4.4 says that $f(x) = f(a) + \int_a^x Df(t) dt$ whenever it makes sense to write $\int_a^x Df(t) dt$, if Df is interpreted as a weak derivative. Somehow, this is the way it ought to be. It follows from the fundamental theorem of calculus that $f'(x)$ exists for a.e. x in the classical sense where the derivative is taken in the sense of a limit of difference quotients and $f'(x) = Df(x)$. This raises an interesting question. Suppose f is continuous on $[a, b]$ and $f'(x)$ exists in the classical sense for a.e. x . Does it follow that $f(x) = f(a) + \int_a^x f'(t) dt$? The answer is no. You can build such an example from the Cantor function which is increasing and has a derivative a.e. which equals 0 a.e. and yet climbs from 0 to 1, Problem 4 on Page 268. Thus, in a sense weak derivatives are more agreeable than the classical ones.

16.5 Definition of Hausdorff Measures

First I will discuss some outer measures. In all that is done here, $\alpha(p)$ will be the volume of the ball in \mathbb{R}^p which has radius 1. Hausdorff measures are very geometrically motivated and so the norm in \mathbb{R}^p will be the usual Euclidean norm unless indicated otherwise.

Definition 16.5.1 *For a set E , denote by $r(E)$ the number which is half the diameter of E . Thus $r(E) \equiv \frac{1}{2} \sup \{ |x - y| : x, y \in E \} \equiv \frac{1}{2} \text{diam}(E)$. Let $E \subseteq \mathbb{R}^p$. $\mathcal{H}_\delta^s(E) \equiv \inf \{ \sum_{j=1}^\infty \beta(s)(r(C_j))^s : E \subseteq \cup_{j=1}^\infty C_j, r(C_j) < \delta \}$. Define $\mathcal{H}^s(E) \equiv \lim_{\delta \rightarrow 0+} \mathcal{H}_\delta^s(E)$. In case $s = 0$ and E is an infinite set, then for small enough δ , $\mathcal{H}_\delta^0(E) > m$ for any positive m so $\mathcal{H}^0(E) = \infty$ and in case E is a finite set, then $\mathcal{H}_\delta^0(E)$ will clearly be the number of things in E for all δ small enough. Thus \mathcal{H}^0 is just counting measure and the \mathcal{H}_δ^0 are outer measures converging to \mathcal{H}^0 provided we define $\mathcal{H}^0(\emptyset) \equiv 0$.*

Note that $\mathcal{H}_\delta^s(E)$ if you make δ smaller, $\mathcal{H}_\delta^s(E)$ will become larger and so the limit clearly exists.

In the above definition, $\beta(s)$ is an appropriate **positive constant** depending on s . It will turn out that for p an integer, $\beta(p) = \alpha(p)$ where $\alpha(p)$ is the Lebesgue measure of the unit ball, $B(0, 1)$ where the Euclidean norm is used to determine this ball.

Lemma 16.5.2 \mathcal{H}^s and \mathcal{H}_δ^s are outer measures for all $s \geq 0$.

Proof: The case $s = 0$ comes directly from the definition so assume $s > 0$. If $A \subseteq B$, then $\mathcal{H}^s(A) \leq \mathcal{H}^s(B)$ with similar assertions valid for \mathcal{H}_δ^s . To see that $\mathcal{H}_\delta^s(\emptyset) = 0$, let $C_j \equiv B(0, \varepsilon^{1/s} 2^{-(j+1)/s})$ where $\varepsilon^{1/s} < \delta$ so that $\mathcal{H}_\delta^s(\emptyset) \leq \sum_{j=1}^\infty \beta(s) \left(\varepsilon^{1/s} 2^{-(j+1)/s} \right)^s = \beta(s) \frac{\varepsilon}{2}$. Since ε is arbitrary, $\mathcal{H}^s(\emptyset) = 0$.

Suppose $E = \cup_{i=1}^\infty E_i$ and $\mathcal{H}_\delta^s(E_i) < \infty$ for each i . Let $\{C_j^i\}_{j=1}^\infty$ be a covering of E_i with $\sum_{j=1}^\infty \beta(s)(r(C_j^i))^s - \varepsilon/2^i < \mathcal{H}_\delta^s(E_i)$ and $\text{diam}(C_j^i) \leq \delta$. Then

$$\mathcal{H}_\delta^s(E) \leq \sum_{i=1}^\infty \sum_{j=1}^\infty \beta(s)(r(C_j^i))^s \leq \sum_{i=1}^\infty \mathcal{H}_\delta^s(E_i) + \varepsilon/2^i \leq \varepsilon + \sum_{i=1}^\infty \mathcal{H}_\delta^s(E_i).$$

It follows that since $\varepsilon > 0$ is arbitrary, $\mathcal{H}_\delta^s(E) \leq \sum_{i=1}^\infty \mathcal{H}_\delta^s(E_i)$ which shows \mathcal{H}_δ^s is an outer measure. Now notice that $\mathcal{H}_\delta^s(E)$ is increasing as $\delta \rightarrow 0$. Picking a sequence δ_k decreasing to 0, the monotone convergence theorem implies $\mathcal{H}^s(E) \leq \sum_{i=1}^\infty \mathcal{H}^s(E_i)$. ■

The outer measure \mathcal{H}^s is called s dimensional Hausdorff measure when restricted to the σ algebra of \mathcal{H}^s measurable sets. It is automatically a complete measure meaning that if $E \subseteq F$ where $\mathcal{H}^s(F) = 0$ then E is measurable. This follows from Theorem 9.5.4.

Next I will show the σ algebra of \mathcal{H}^s measurable sets includes the Borel sets.

16.6 Properties of Hausdorff Measure

Using Theorem 9.6.1 on Page 248, the following is obtained.

Theorem 16.6.1 The σ algebra of \mathcal{H}^s measurable sets contains the Borel sets and \mathcal{H}^s has the property that for all $E \subseteq \mathbb{R}^p$, there exists a Borel set $F \supseteq E$ such that $\mathcal{H}^s(F) = \mathcal{H}^s(E)$.

Proof: Let $\text{dist}(A, B) = 2\delta_0 > 0$. Is it the case that $\mathcal{H}^s(A) + \mathcal{H}^s(B) = \mathcal{H}^s(A \cup B)$? This is what is needed to use Theorem 9.6.1 about measurable sets including the Borel sets.

Let $\{C_j\}_{j=1}^\infty$ be a covering of $A \cup B$ such that $\text{diam}(C_j) \leq \delta < \delta_0$ for each j and

$$\mathcal{H}_\delta^s(A \cup B) + \varepsilon > \sum_{j=1}^\infty \beta(s)(r(C_j))^s.$$

Thus $\mathcal{H}_\delta^s(A \cup B) + \varepsilon > \sum_{j \in J_1} \beta(s)(r(C_j))^s + \sum_{j \in J_2} \beta(s)(r(C_j))^s$ where

$$J_1 = \{j : C_j \cap A \neq \emptyset\}, J_2 = \{j : C_j \cap B \neq \emptyset\}.$$

Recall $\text{dist}(A, B) = 2\delta_0$, $J_1 \cap J_2 = \emptyset$. It follows $\mathcal{H}_\delta^s(A \cup B) + \varepsilon > \mathcal{H}_\delta^s(A) + \mathcal{H}_\delta^s(B)$. Letting $\delta \rightarrow 0$, and noting $\varepsilon > 0$ was arbitrary, yields $\mathcal{H}^s(A \cup B) \geq \mathcal{H}^s(A) + \mathcal{H}^s(B)$. Equality holds because \mathcal{H}^s is an outer measure. By Theorem 9.6.1, \mathcal{H}^s is a Borel measure.

To verify the second assertion, note there is no loss of generality in letting $\mathcal{H}^s(E) < \infty$. Let $E \subseteq \cup_{j=1}^\infty C_j$, $r(C_j) < \delta$, and

$$\mathcal{H}_\delta^s(E) + \delta > \sum_{j=1}^\infty \beta(s)(r(C_j))^s. \quad (16.10)$$

Let $F_\delta = \bigcup_{j=1}^\infty \overline{C_j}$. Thus $F_\delta \supseteq E$ and

$$\mathcal{H}_\delta^s(E) \leq \mathcal{H}_\delta^s(F_\delta) \leq \sum_{j=1}^\infty \beta(s)(r(\overline{C_j}))^s = \sum_{j=1}^\infty \beta(s)(r(C_j))^s < \delta + \mathcal{H}_\delta^s(E).$$

Let $\delta_k \rightarrow 0$ and let $F = \bigcap_{k=1}^\infty F_{\delta_k}$. Then $F \supseteq E$ and

$$\mathcal{H}_{\delta_k}^s(E) \leq \mathcal{H}_{\delta_k}^s(F) \leq \mathcal{H}_{\delta_k}^s(F_{\delta_k}) \leq \delta_k + \mathcal{H}_{\delta_k}^s(E).$$

Letting $k \rightarrow \infty$, $\mathcal{H}^s(E) \leq \mathcal{H}^s(F) \leq \mathcal{H}^s(E)$

We can also arrange to have F containing E be a G_δ set. In 16.10, replace C_j with $C_j + B(0, \eta_j)$ which is an open set having diameter no more than $\text{diam}(C_j) + 2\eta_j$ so by taking η_j small enough, we can replace each C_j with an open set O_j in such a way as to preserve 16.10 with C_j replaced with O_j and also $r(O_j) < \delta$. Then letting $V_\delta \equiv \bigcup_j O_j$,

$$\mathcal{H}_\delta^s(E) \leq \mathcal{H}_\delta^s(V_\delta) \leq \sum_{j=1}^\infty \beta(s)(r(O_j))^s < \delta + \mathcal{H}_\delta^s(E).$$

Then let $G = \bigcap_k V_{\delta_k}$ where $\delta_k \rightarrow 0$ and let the V_{δ_k} be decreasing as k increases, each V_δ containing E . Then for each δ , $\mathcal{H}_{\delta_k}^s(E) \leq \mathcal{H}_{\delta_k}^s(G) < \delta + \mathcal{H}_{\delta_k}^s(E)$. Let $k \rightarrow \infty$ to find $\mathcal{H}^s(E) \leq \mathcal{H}^s(G) \leq \mathcal{H}^s(E)$ as before. ■

A measure satisfying the conclusion of Theorem 16.6.1 is called a Borel regular measure.

16.7 \mathcal{H}^p and m_p

Next I will compare \mathcal{H}^p and m_p . First recall this covering theorem which is a summary of Corollary 9.12.5 found on Page 265.

Theorem 16.7.1 *Let $E \subseteq \mathbb{R}^p$ and let \mathcal{F} be a collection of balls of bounded radii such that \mathcal{F} covers E in the sense of Vitali. Then there exists a countable collection of disjoint balls from \mathcal{F} , $\{B_j\}_{j=1}^\infty$, such that $\overline{m_p}(E \setminus \bigcup_{j=1}^\infty B_j) = 0$.*

Recall the following interesting lemma stated here for convenience. It is Lemma 11.7.2.

Lemma 16.7.2 *Every open set U in \mathbb{R}^p is a countable disjoint union of half open boxes of the form $Q \equiv \prod_{i=1}^p [a_i, a_i + 2^{-k})$ where $a_i = l2^{-k}$ for l some integer.*

Lemma 16.7.3 *If $S \subseteq \mathbb{R}^p$ and $m_p(S) = 0$, then $\mathcal{H}^p(S) = \mathcal{H}_\delta^p(S) = 0$. Also, there exists a constant k such that $\mathcal{H}^p(E) \leq km_p(E)$ for all E Borel $k \equiv \frac{\beta(p)}{\alpha(p)}$. Also, if $Q_0 \equiv [0, 1)^p$, the unit cube, then $\infty > \mathcal{H}^p([0, 1)^p) > 0$.*

Proof: Suppose first $m_p(S) = 0$. Without loss of generality, S is bounded. Then by outer regularity, there exists a bounded open V containing S and $m_p(V) < \varepsilon$. For each $x \in S$, there exists a ball B_x such that $\widehat{B_x} \subseteq V$ and $\delta > r(\widehat{B_x})$. By the Vitali covering theorem there is a sequence of disjoint balls $\{B_k\}$ such that $\{\widehat{B_k}\}$ covers S . Here $\widehat{B_k}$ has

the same center as B_k but 5 times the radius. Then letting $\alpha(p)$ be the Lebesgue measure of the unit ball in \mathbb{R}^p

$$\mathcal{H}_\delta^p(S) \leq \sum_k \beta(p) r(\widehat{B_k})^p = \frac{\beta(p)}{\alpha(p)} 5^p \sum_k \alpha(p) r(B_k)^p \leq \frac{\beta(p)}{\alpha(p)} 5^p m_p(V) < \frac{\beta(p)}{\alpha(p)} 5^p \varepsilon$$

Since ε is arbitrary, this shows $\mathcal{H}_\delta^p(S) = 0$ and $\mathcal{H}^p(S) \equiv \lim_{\delta \rightarrow 0} \mathcal{H}_\delta^p(S) = 0$.

Letting U be an open set and $\delta > 0$, consider all balls B contained in U which have diameters less than δ . This is a Vitali covering of U and therefore by Theorem 16.7.1, there exists $\{B_i\}$, a sequence of disjoint balls of radii less than δ contained in U such that $\cup_{i=1}^\infty B_i$ differs from U by a set of Lebesgue measure zero. Let $\alpha(p)$ be the Lebesgue measure of the unit ball in \mathbb{R}^p . Then from what was just shown,

$$\begin{aligned} \mathcal{H}_\delta^p(U) &= \mathcal{H}_\delta^p(\cup_i B_i) \leq \sum_{i=1}^\infty \beta(p) r(B_i)^p = \frac{\beta(p)}{\alpha(p)} \sum_{i=1}^\infty \alpha(p) r(B_i)^p \\ &= \frac{\beta(p)}{\alpha(p)} \sum_{i=1}^\infty m_p(B_i) = \frac{\beta(p)}{\alpha(p)} m_p(U) \equiv k m_p(U), \quad k \equiv \frac{\beta(p)}{\alpha(p)} \end{aligned}$$

Now letting E be Lebesgue measurable, it follows from the outer regularity of m_p there exists a decreasing sequence of open sets, $\{V_i\}$ containing E such that $m_p(V_i) \rightarrow m_p(E)$. Then from the above, $\mathcal{H}_\delta^p(E) \leq \lim_{i \rightarrow \infty} \mathcal{H}_\delta^p(V_i) \leq \lim_{i \rightarrow \infty} k m_p(V_i) = k m_p(E)$. Since $\delta > 0$ is arbitrary, it follows that also $\mathcal{H}^p(E) \leq k m_p(E)$. This proves the first part of the lemma and that $\infty > \mathcal{H}^p([0, 1]^p)$.

If $\mathcal{H}^p([0, 1]^p) = 0$, it follows $\mathcal{H}^p(\mathbb{R}^p) = 0$ because \mathbb{R}^p is the countable union of translates of $Q_0 \equiv [0, 1]^p$ and it is clear that \mathcal{H}^p is translation invariant. Since each \mathcal{H}_δ^p is no larger than \mathcal{H}^p , $\mathcal{H}_\delta^p(\mathbb{R}^p) = 0$. Therefore, there exists a sequence of sets, $\{C_i\}$ each having diameter less than δ such that the union of these sets equals \mathbb{R}^p but $1 > \sum_{i=1}^\infty \beta(p) r(C_i)^p$. Now let B_i be a ball having radius r_i equal to $\text{diam}(C_i) = 2r(C_i)$ which contains C_i . These B_i cover \mathbb{R}^p , $\frac{1}{2}r_i = r(C_i)$. It follows that

$$1 > \sum_{i=1}^\infty \beta(p) r(C_i)^p = \sum_{i=1}^\infty \frac{\beta(p)}{\alpha(p) 2^p} m_p(B_i) = \infty,$$

a contradiction. This shows that $\mathcal{H}^p([0, 1]^p) > 0$. ■

Note that the above shows that $\mathcal{H}^p([-n, n]^p)$ is always a finite positive real number for $n \in \mathbb{N}$.

Theorem 16.7.4 *By choosing $\beta(p)$ properly, one can obtain $\mathcal{H}^p = m_p$ on all Lebesgue measurable sets.*

Proof: Define $l = \frac{m_p(Q_0)}{\mathcal{H}^p(Q_0)}$ where $Q_0 = [0, 1]^p$ is the half open unit cube in \mathbb{R}^p . It follows then that $l = \frac{m_p(Q)}{\mathcal{H}^p(Q)}$ where $Q = \prod_{i=1}^p [a_i, a_i + 2^{-k}]$ where $a_i = l 2^{-k}$ for l some integer because of translation invariance of both measures and that Q_0 is the union of such Q . By Lemma 16.7.2, $l \mathcal{H}^p(V) = m_p(V)$ for any V open. Letting V_n be an increasing sequence of bounded open sets whose union is \mathbb{R}^p , it follows that the set of Borel E satisfying $l \mathcal{H}^p(E \cap V_n) = m_p(E \cap V_n)$ is a σ algebra which contains the open sets and so this equation is true for all Borel sets. Letting $n \rightarrow \infty$, $l \mathcal{H}^p(F) = m_p(F)$ for any Borel F . For E Lebesgue measurable, there is F Borel contained in E with $m_p(E \setminus F) = 0$ and so

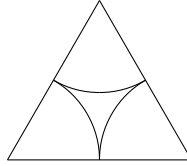
F is \mathcal{H}^p measurable as is $E \setminus F$ because $\mathcal{H}^p(E \setminus F) = 0$ from Lemma 16.7.3 and \mathcal{H}^p is an outer measure. Recall that sets having outer measure 0 end up being measurable sets. Thus $E = F \cup (E \setminus F)$ is \mathcal{H}^p measurable also. This implies from Theorem 16.6.1 and Proposition 11.1.2 that if E is Lebesgue measurable, there is Borel $F \supseteq E$ such that $m_p(E) = m_p(F) = l\mathcal{H}^p(F) = l\mathcal{H}^p(E)$. Now choose $\beta(p)$ to make the constant $l = 1$. ■

The exact determination of $\beta(p)$ is more technical.

16.8 Technical Considerations

Let $\alpha(p)$ be the volume of the unit ball in \mathbb{R}^p . Thus the volume of $B(\mathbf{0}, r)$ in \mathbb{R}^p is $\alpha(p)r^p$ from the change of variables formula. There is a very important and interesting inequality known as the isodiametric inequality which says that if A is any set in \mathbb{R}^p , then

$$\bar{m}_p(A) \leq \alpha(p)(2^{-1}\text{diam}(A))^p = \alpha(p)r(A)^p.$$



This inequality may seem obvious at first but it is not really. The reason it is not is that there are sets which are not subsets of any sphere having the same diameter as the set. For example, consider an equilateral triangle. You have to include the vertices and so the center of such a ball would need to be closer to each vertex than the radius of the small circles. See the above picture

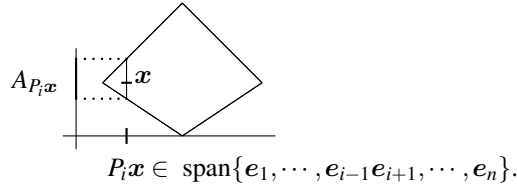
Lemma 16.8.1 Let $f : \mathbb{R}^{p-1} \rightarrow [0, \infty)$ be Borel measurable and let

$$S = \{(x, y) : |y| < f(x)\}.$$

Then S is a Borel set in \mathbb{R}^p .

Proof: Set s_k be an increasing sequence of Borel measurable functions converging pointwise to f . $s_k(x) = \sum_{m=1}^{N_k} c_m^k \chi_{E_m^k}(x)$. Let $S_k = \bigcup_{m=1}^{N_k} E_m^k \times (-c_m^k, c_m^k)$. Then $(x, y) \in S_k$ if and only if $f(x) > 0$ and $|y| < s_k(x) \leq f(x)$. It follows that $S_k \subseteq S_{k+1}$ and $S = \bigcup_{k=1}^{\infty} S_k$. But each S_k is a Borel set and so S is also a Borel set. ■

Let P_i be the projection onto $\text{span}(e_1, \dots, e_{i-1}, e_{i+1}, \dots, e_p)$ where the e_k are the standard basis vectors in \mathbb{R}^p , e_k being the vector having a 1 in the k^{th} slot and a 0 elsewhere. Thus $P_i x \equiv \sum_{j \neq i} x_j e_j$. Also let $A_{P_i x} \equiv \{x_i : (x_1, \dots, x_i, \dots, x_p) \in A\}$



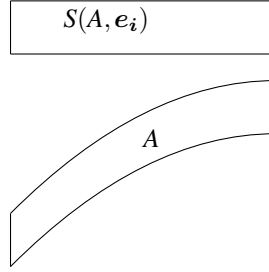
Lemma 16.8.2 Let $A \subseteq \mathbb{R}^p$ be a Borel set. Then $P_i x \rightarrow m(A_{P_i x})$ is a Borel measurable function defined on $P_i(\mathbb{R}^p)$.

Proof: From Proposition 10.14.4, $A_{P_i x}$ is measurable if S is product measurable. By Theorem 10.14.9, the Borel sets are product measurable and $P_i x \rightarrow m(A_{P_i x})$ is measurable, in fact product measurable. Therefore, the desired conclusion of this lemma follows. ■

16.8.1 Steiner Symmetrization

Definition 16.8.3 Define $S(A, e_i) \equiv \{x \equiv P_i x + y e_i : |y| < 2^{-1} m(A_{P_i x})\}$.

Here is a picture of the idea used in producing $S(A, e_i)$ from A . The one on the top is $S(A, e_i)$. The two sets have the same area but $S(A, e_i)$ smaller diameter than A .



Lemma 16.8.4 Let A be a Borel subset of \mathbb{R}^p . Then $S(A, e_i)$ satisfies

$$P_i x + y e_i \in S(A, e_i) \text{ if and only if } P_i x - y e_i \in S(A, e_i),$$

$$S(A, e_i) \text{ is a Borel set in } \mathbb{R}^p,$$

$$m_p(S(A, e_i)) = m_p(A), \quad (16.11)$$

$$\text{diam}(S(A, e_i)) \leq \text{diam}(A). \quad (16.12)$$

Proof: The first assertion is obvious from the definition. The Borel measurability of $S(A, e_i)$ follows from the definition and Lemmas 16.8.2 and 16.8.1. To show 16.11,

$$\begin{aligned} m_p(S(A, e_i)) &= \int_{P_i \mathbb{R}^p} \int_{-2^{-1} m(A_{P_i x})}^{2^{-1} m(A_{P_i x})} dx_i dx_1 \cdots dx_{i-1} dx_{i+1} \cdots dx_p = \\ &= \int_{P_i \mathbb{R}^p} m(A_{P_i x}) dx_1 \cdots dx_{i-1} dx_{i+1} \cdots dx_p = \int_{P_i \mathbb{R}^p} \int_{\mathbb{R}} \mathcal{X}_A dx_i dx_1 \cdots dx_{i-1} dx_{i+1} \cdots dx_p = m_p(A) \end{aligned}$$

Now suppose x_1 and $x_2 \in S(A, e_i)$, and $x_1 = P_i x_1 + y_1 e_i$, $x_2 = P_i x_2 + y_2 e_i$.

Then $y_1 \in \left[-\frac{m(A_{P_i x_1})}{2}, \frac{m(A_{P_i x_1})}{2} \right]$, $y_2 \in \left[-\frac{m(A_{P_i x_2})}{2}, \frac{m(A_{P_i x_2})}{2} \right]$. There exists

$$x_{1i} \in [\inf A_{P_i x_1}, \sup A_{P_i x_1}] \text{ and } x_{2i} \in [\inf A_{P_i x_2}, \sup A_{P_i x_2}]$$

such that $x_{1i} \in A_{P_i x_1}$ and $x_{2i} \in A_{P_i x_2}$. The second pair of intervals is at least as long as the corresponding interval in the first pair and the second pair are not necessarily centered at the same point. Therefore, such an x_{1i} and x_{2i} can be chosen such that $|x_{2i} - x_{1i}| \geq |y_1 - y_2|$ and so $\hat{x}_1 \equiv P_i x_1 + x_{1i} e_i$ and $\hat{x}_2 \equiv P_i x_2 + x_{2i} e_i$ are in A and $|\hat{x}_1 - \hat{x}_2| \geq |x_1 - x_2|$ so the diameter of $S(A, e_i)$ is no more than the diameter of A as claimed. ■

The next lemma says that if A is already symmetric with respect to the j^{th} direction, then this symmetry is not destroyed by taking $S(A, e_i)$.

Lemma 16.8.5 Suppose A is a Borel set in \mathbb{R}^p such that $P_j x + e_j x_j \in A$ if and only if $P_j x + (-x_j) e_j \in A$. Then if $i \neq j$, $P_j x + e_j x_j \in S(A, e_i)$ if and only if $P_j x + (-x_j) e_j \in S(A, e_i)$.

Proof: By definition, $P_j \mathbf{x} + e_j x_j \in S(A, e_i)$ if and only if $|x_i| < 2^{-1} m(A_{P_i(P_j \mathbf{x} + e_j x_j)})$. Now $x_i \in A_{P_i(P_j \mathbf{x} + e_j x_j)}$ if and only if $x_i \in A_{P_i(P_j \mathbf{x} + (-x_j) e_j)}$ by the assumption on A which says that A is symmetric in the e_j direction. Hence $P_j \mathbf{x} + e_j x_j \in S(A, e_i)$ if and only if $|x_i| < 2^{-1} m(A_{P_i(P_j \mathbf{x} + (-x_j) e_j)})$ if and only if $P_j \mathbf{x} + (-x_j) e_j \in S(A, e_i)$. ■

16.8.2 The Isodiametric Inequality

The next theorem is called the isodiametric inequality. It is the key result used to compare Lebesgue and Hausdorff measures.

Theorem 16.8.6 *Let A be any Lebesgue measurable set in \mathbb{R}^p . Then it follows that $m_p(A) \leq \alpha(p)(r(A))^p$.*

Proof: Suppose first that A is Borel. Let $A_1 = S(A, e_1)$ and $A_k = S(A_{k-1}, e_k)$. Then by Lemma 16.8.4, A_p is a Borel set, $\text{diam}(A_p) \leq \text{diam}(A)$, $m_p(A_p) = m_p(A)$ and A_p is symmetric. Thus $\mathbf{x} \in A_p$ if and only if $-\mathbf{x} \in A_p$. It follows that $A_p \subseteq \overline{B(\mathbf{0}, r(A_p))}$. If $\mathbf{x} \in A_p \setminus \overline{B(\mathbf{0}, r(A_p))}$, then $-\mathbf{x} \in A_p \setminus \overline{B(\mathbf{0}, r(A_p))}$ and so $\text{diam}(A_p) \geq 2|\mathbf{x}| > \text{diam}(A_p)$. Therefore, there is no such \mathbf{x} and $m_p(A_p) \leq \alpha(p)(r(A_p))^p \leq \alpha(p)(r(A))^p$. It remains to establish this inequality for arbitrary measurable sets. Letting A be such a set, let $\{K_k\}$ be an increasing sequence of compact subsets of A such that $m_p(A) = \lim_{k \rightarrow \infty} m_p(K_k)$. Then

$$m_p(A) = \lim_{k \rightarrow \infty} m_p(K_k) \leq \limsup_{k \rightarrow \infty} \alpha(p)(r(K_k))^p \leq \alpha(p)(r(A))^p. \quad \blacksquare$$

16.9 The Proper Value of $\beta(p)$

I will show that the proper determination of $\beta(p)$ is $\alpha(p)$, the volume of the unit ball. Since $\beta(p)$ has been adjusted such that $l = 1$ in Theorem 16.7.4, $m_p(B(\mathbf{0}, 1)) = \mathcal{H}^p(B(\mathbf{0}, 1))$. There exists a covering of $B(\mathbf{0}, 1)$ of sets of radii less than δ , $\{C_i\}_{i=1}^\infty$ such that

$$\mathcal{H}_\delta^p(B(\mathbf{0}, 1)) + \varepsilon > \sum_i \beta(p) r(C_i)^p$$

Then by Theorem 16.8.6, the isodiametric inequality,

$$\begin{aligned} \mathcal{H}_\delta^p(B(\mathbf{0}, 1)) + \varepsilon &> \sum_i \beta(p) r(C_i)^p = \frac{\beta(p)}{\alpha(p)} \sum_i \alpha(p) r(\overline{C_i})^p \\ &\geq \frac{\beta(p)}{\alpha(p)} \sum_i m_p(\overline{C_i}) \geq \frac{\beta(p)}{\alpha(p)} m_p(B(\mathbf{0}, 1)) = \frac{\beta(p)}{\alpha(p)} \mathcal{H}^p(B(\mathbf{0}, 1)) \end{aligned}$$

Now taking the limit as $\delta \rightarrow 0$, $\mathcal{H}^p(B(\mathbf{0}, 1)) + \varepsilon \geq \frac{\beta(p)}{\alpha(p)} \mathcal{H}^p(B(\mathbf{0}, 1))$ and since $\varepsilon > 0$ is arbitrary, this shows $\alpha(p) \geq \beta(p)$.

By the Vitali covering theorem in Corollary 9.12.5, there exists a sequence of disjoint balls of radius no more than δ , $\{B_i\}$ such that $B(\mathbf{0}, 1) = (\cup_{i=1}^\infty B_i) \cup N$, where $m_p(N) = 0$. Then $\mathcal{H}_\delta^p(N) = 0$ can be concluded because $\mathcal{H}_\delta^p \leq \mathcal{H}^p$ and Lemma 16.7.3. Using $m_p(B(\mathbf{0}, 1)) = \mathcal{H}^p(B(\mathbf{0}, 1))$ again,

$$\mathcal{H}_\delta^p(B(\mathbf{0}, 1)) = \mathcal{H}_\delta^p(\cup_i B_i) \leq \sum_{i=1}^\infty \beta(p) r(B_i)^p = \frac{\beta(p)}{\alpha(p)} \sum_{i=1}^\infty \alpha(p) r(B_i)^p$$

$$= \frac{\beta(p)}{\alpha(p)} \sum_{i=1}^{\infty} m_p(B_i) = \frac{\beta(p)}{\alpha(p)} m_p(\cup_i B_i) = \frac{\beta(p)}{\alpha(p)} m_p(B(\mathbf{0}, 1)) = \frac{\beta(p)}{\alpha(p)} \mathcal{H}^p(B(\mathbf{0}, 1))$$

which implies $\alpha(p) \leq \beta(p)$ and so the two are equal. This proves that if $\alpha(p) = \beta(p)$, then the $\mathcal{H}^p = m_p$ on the measurable sets of \mathbb{R}^p .

This gives another way to think of Lebesgue measure which is a particularly nice way because it is coordinate free, depending only on the notion of distance.

For $s < p$, note that \mathcal{H}^s is not a Radon measure because it will not generally be finite on compact sets. For example, let $p = 2$ and consider $\mathcal{H}^1(L)$ where L is a line segment joining $(0, 0)$ to $(1, 0)$. Then $\mathcal{H}^1(L)$ is no smaller than $\mathcal{H}^1(L)$ when L is considered a subset of \mathbb{R}^1 , $p = 1$. Thus by what was just shown, $\mathcal{H}^1(L) \geq 1$. Hence $\mathcal{H}^1([0, 1] \times [0, 1]) = \infty$. The situation is this: L is a one-dimensional object inside \mathbb{R}^2 and \mathcal{H}^1 is giving a one-dimensional measure of this object. In fact, Hausdorff measures can make such heuristic remarks as these precise. Define the Hausdorff dimension of a set A , as

$$\dim(A) = \inf\{s : \mathcal{H}^s(A) = 0\}$$

16.10 A Formula for $\alpha(p)$

What is $\alpha(p)$ for p a positive integer? Let p be a positive integer. Theorem 14.4.1 on Page 405 says that

Theorem 16.10.1 $\alpha(p) = \pi^{p/2}(\Gamma(p/2 + 1))^{-1}$ where $\Gamma(s)$ is the gamma function $\Gamma(s) = \int_0^\infty e^{-t} t^{s-1} dt$.

From now on, in the definition of Hausdorff measure, it will always be the case that $\beta(s) = \alpha(s)$. As shown above, this is the right thing to have $\beta(s)$ to equal if s is a positive integer because this yields the important result that Hausdorff measure is the same as Lebesgue measure. Note the formula, $\pi^{s/2}(\Gamma(s/2 + 1))^{-1}$ makes sense for any $s \geq 0$.

Chapter 17

The Area Formula

I am grateful to those who have found errors in this material, some of which were egregious. I would not have found these mistakes because I never teach this material and I don't use it in my research. I do think it is wonderful mathematics however.

17.1 Estimates for Hausdorff Measure

This section is on estimates which relate Hausdorff measure to Lebesgue measure. This will allow a geometric motivation for measures on Lipschitz manifolds.

The main case will be for h a Lipschitz function, $|h(x) - h(y)| \leq K|x - y|$ defined on \mathbb{R}^n . This is no loss of generality because of Theorem 16.3.6. However, the main presentation will include more general situations than this. One uses the differentiability of h off a set of measure zero to show the existence of disjoint Borel sets E on which h is Lipschitz with its inverse also being Lipschitz on $h(E)$.

The following lemma states that Lipschitz maps take sets of measure zero to sets of measure zero. It also gives a convenient estimate. This involves the consideration of \mathcal{H}^n as an outer measure. Thus it is not necessary to know the set B is measurable.

In fact, one only needs to have h locally Lipschitz in much of what follows.

Definition 17.1.1 Let $h : \mathbb{R}^n \rightarrow \mathbb{R}^m$. This function is said to be locally Lipschitz if for every $x \in \mathbb{R}^n$, there exists a ball B_x containing x and a constant K_x such that for all $y, z \in B_x$,

$$|h(z) - h(y)| \leq K_x |z - y|$$

Lemma 17.1.2 If h is Lipschitz with Lipschitz constant K then for $B \subseteq \mathbb{R}^n$,

$$\mathcal{H}^n(h(B)) \leq K^n \mathcal{H}^n(B)$$

Also, if T is a set in \mathbb{R}^n , $m_n(T) = 0$, then $\mathcal{H}^n(h(T)) = 0$. It is not necessary that h be one to one.

Proof: If $\mathcal{H}^n(B) = \infty$, there is nothing to show. Assume $\mathcal{H}^n(B) < \infty$. Let $\{C_i\}_{i=1}^\infty$ cover B with each having diameter less than δ and let this cover be such that

$$\sum_i \beta(n) \frac{1}{2} \text{diam}(C_i)^n < \mathcal{H}_\delta^n(B) + \epsilon$$

Then $\{h(C_i)\}$ covers $h(B)$ and each set has diameter no more than $K\delta$. Then

$$\begin{aligned} \mathcal{H}_{K\delta}^n(h(B)) &\leq \sum_i \beta(n) \left(\frac{1}{2} \text{diam}(h(C_i)) \right)^n \\ &\leq K^n \sum_i \beta(n) \left(\frac{1}{2} \text{diam}(C_i) \right)^n \leq K^n (\mathcal{H}_\delta^n(B) + \epsilon) \end{aligned}$$

Since ϵ is arbitrary, this shows that $\mathcal{H}_{K\delta}^n(h(B)) \leq K^n \mathcal{H}_\delta^n(B)$. Now take a limit as $\delta \rightarrow 0$. The second claim follows from $m_n = \mathcal{H}^n$ on Lebesgue measurable sets of \mathbb{R}^n . ■

Lemma 17.1.3 If h is locally Lipschitz and $m_n(T) = 0$, then $\mathcal{H}^n(h(T)) = 0$. It is not necessary that h be one to one.

Proof: Let $T_k \equiv \{x \in T : h \text{ has Lipschitz constant } k \text{ near } x\}$. Thus $T = \bigcup_k T_k$. I will show $h(T_k)$ has \mathcal{H}^n measure zero and then it will follow that $h(T) = \bigcup_{k=1}^\infty h(T_k)$, the $h(T_k)$ increasing in k , must also have measure zero.

Let $\varepsilon > 0$ be given. By outer regularity, there exists an open set V containing T_k such that $m_n(V) < \varepsilon$. For $x \in T_k$ it follows there exists $r_x < 1$ such that the ball centered at x with radius r_x is contained in V and in this ball, h has Lipschitz constant k . By the Besicovitch covering theorem, Theorem 4.5.8, there are N_n sets of these balls $\{\mathcal{G}_1, \dots, \mathcal{G}_{N_n}\}$ such that the balls in \mathcal{G}_k are disjoint and the union of all balls in the N_n sets covers T_k . Then

$$\begin{aligned} \mathcal{H}^n(h(T_k)) &\leq \sum_{k=1}^{N_n} \sum \{\mathcal{H}^n(h(B)) : B \in \mathcal{G}_k\} \\ &\leq N_n k^n m_n(V) < N_n k^n \varepsilon \end{aligned}$$

Since ε is arbitrary, this shows that $\mathcal{H}^n(h(T_k)) = 0$. Hence $\mathcal{H}^n(h(T)) = 0$ also, since it is the limit of the $\mathcal{H}^n(h(T_k))$. ■

Lemma 17.1.4 *If S is a Lebesgue measurable set in \mathbb{R}^n and h is Lipschitz or locally Lipschitz then $h(S)$ is \mathcal{H}^n measurable. Also, if h is Lipschitz with constant K , $\mathcal{H}^n(h(S)) \leq K^n m_n(S)$. It is not necessary that h be one to one.*

Proof: The estimate follows from Lemma 17.1.2 or 17.1.3 and the observation that, as shown before, Theorem 16.7.4, if S is Lebesgue measurable in \mathbb{R}^n , then $\mathcal{H}^n(S) = m_n(S)$. The estimate also shows that h maps sets of Lebesgue measure zero to sets of \mathcal{H}^n measure zero. Why is $h(S)$ \mathcal{H}^n measurable if S is Lebesgue measurable? This follows from completeness of \mathcal{H}^n . Indeed, let F be F_σ and contained in S with $m_n(S \setminus F) = 0$. Then $h(S) = h(S \setminus F) \cup h(F)$. The second set is Borel and the first has \mathcal{H}^n measure zero. By completeness of \mathcal{H}^n , $h(S)$ is \mathcal{H}^n measurable. ■

Recall Theorem 1.5.5 on Page 23. This is stated here for convenience.

Theorem 17.1.5 *Let F be an $m \times p$ matrix where $m \geq p$. Then there exists an $m \times p$ matrix R and a $p \times p$ matrix U such that*

$$F = RU, \quad U = U^*,$$

*all eigenvalues of U are non negative, $U^2 = F^*F$, $R^*R = I$, and $|Rx| = |x|$.*

Thus, if $h : \mathbb{R}^p \rightarrow \mathbb{R}^m, m \geq p$, and $Dh(x)$ exists, then $Dh(x) = R(x)U(x)$ where

$$(U(x)u, v) = (U(x)v, u), (U(x)u, u) \geq 0$$

and $R^*R = I$ so R preserves lengths. Recall that R^* is the adjoint defined by $(Rx, y) = (x, R^*y)$. This convention will be used in what follows.

Lemma 17.1.6 *In this situation where $R^*R = I$, $|R^*u| \leq |u|$.*

Proof: First note that $(u - RR^*u, RR^*u) = (u, RR^*u) - |RR^*u|^2 = |R^*u|^2 - |R^*u|^2 = 0$, and so

$$|u|^2 = |u - RR^*u + RR^*u|^2 = |u - RR^*u|^2 + |RR^*u|^2 = |u - RR^*u|^2 + |R^*u|^2. \quad \blacksquare$$

Then the following corollary follows from Lemma 17.1.6.

Corollary 17.1.7 *Let $T \subseteq \mathbb{R}^m$. Then $\mathcal{H}^n(T) \geq \mathcal{H}^n(R^*T)$.*

Hausdorff measure makes possible a unified development of p dimensional area. As in the case of Lebesgue measure, the first step in this is to understand basic considerations related to linear transformations.

Lemma 17.1.8 *Let $R \in \mathcal{L}(\mathbb{R}^p, \mathbb{R}^m)$, $p \leq m$, and $R^*R = I$. Then if $A \subseteq \mathbb{R}^p$, $\mathcal{H}^p(RA) = \mathcal{H}^p(A)$. In fact, if $P: \mathbb{R}^p \rightarrow \mathbb{R}^m$ satisfies $|Px - Py| = |x - y|$, then $\mathcal{H}^p(PA) = \mathcal{H}^p(A)$.*

Proof: Now let P be an arbitrary mapping which preserves lengths, like R , and let A be bounded so $P(A)$ is also bounded. Then $P(A) \subseteq \bigcup_{j=1}^{\infty} C_j$, $r(C_j) < \delta$, and $\mathcal{H}_{\delta}^p(PA) + \varepsilon > \sum_{j=1}^{\infty} \alpha(p)(r(C_j))^p$. Since P preserves lengths, it follows P is one to one and P^{-1} is one to one on $P(\mathbb{R}^p)$ and P^{-1} also preserves lengths on $P(\mathbb{R}^p)$. Replacing each C_j with $C_j \cap (PA)$,

$$\mathcal{H}_{\delta}^p(PA) + \varepsilon > \sum_{j=1}^{\infty} \alpha(p)r(C_j \cap (PA))^p = \sum_{j=1}^{\infty} \alpha(p)r(P^{-1}(C_j \cap (PA)))^p \geq \mathcal{H}_{\delta}^p(A).$$

Thus $\mathcal{H}_{\delta}^p(PA) \geq \mathcal{H}_{\delta}^p(A)$. Similarly $\mathcal{H}_{\delta}^p(P^{-1}(PA)) \geq \mathcal{H}_{\delta}^p(PA)$ so $\mathcal{H}_{\delta}^p(A) \geq \mathcal{H}_{\delta}^p(PA)$. Letting $\delta \rightarrow 0$ yields the desired conclusion in the case where A is bounded. For the general case, let $A_r = A \cap B(0, r)$. Then $\mathcal{H}^p(PA_r) = \mathcal{H}^p(A_r)$. Now let $r \rightarrow \infty$. ■

Lemma 17.1.9 *Let $F \in \mathcal{L}(\mathbb{R}^p, \mathbb{R}^m)$, $p \leq m$, and let $F = RU$ where R and U are described in Theorem 1.5.5 on Page 23. Then if $A \subseteq \mathbb{R}^p$ is Lebesgue measurable,*

$$\mathcal{H}^p(FA) = \det(U)m_p(A).$$

Proof: Using Theorem 11.7.4 on Page 330 and Theorem 16.7.4,

$$\mathcal{H}^p(FA) = \mathcal{H}^p(RUA) = \mathcal{H}^p(UA) = m_p(UA) = \det(U)m_p(A). \quad \blacksquare$$

Definition 17.1.10 *Define J to equal $\det(U)$. Thus*

$$J = \det((F^*F)^{1/2}) = (\det(F^*F))^{1/2}.$$

This is the essential idea for the area formula, but in the area formula, we must consider $\mathbf{h}: \mathbb{R}^p \rightarrow \mathbb{R}^m$ for \mathbf{h} nonlinear and so $\mathbf{h}(U)$ is not a subspace.

17.2 Comparison Theorems

First is a simple lemma which is fairly interesting which involves comparison of two linear transformations. These are Lemmas 5.3.2 and 5.3.1 which follows from fundamental properties of the operator norm. I am stating them here for convenience.

Lemma 17.2.1 *Suppose S, T are linear, defined on a finite dimensional normed linear space, S^{-1} exists, and let $\delta \in (0, 1)$. Then whenever $\|S - T\|$ is small enough, it follows that*

$$\frac{|Tv|}{|Sv|} \in (1 - \delta, 1 + \delta) \quad (17.1)$$

for all $v \neq 0$. Similarly if T^{-1} exists and $\|S - T\|$ is small enough,

$$\frac{|Tv|}{|Sv|} \in (1 - \delta, 1 + \delta)$$

Lemma 17.2.2 *Let S, T be $n \times n$ matrices which are invertible. Then*

$$\mathbf{o}(Tv) = \mathbf{o}(Sv) = \mathbf{o}(v)$$

and if L is a continuous linear transformation such that for $a < b$,

$$\sup_{v \neq 0} \frac{|Lv|}{|Sv|} < b, \quad \inf_{v \neq 0} \frac{|Lv|}{|Sv|} > a$$

If $\|S - T\|$ is small enough, it follows that the same inequalities hold with S replaced with T . Here $\|\cdot\|$ denotes the operator norm.

17.3 The Area Formula

This follows [17] which is where I encountered this material. Let G be an open set and let $\mathbf{h} : G \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^m$ where $m \geq n$ be continuous. Let $D\mathbf{h}(\mathbf{x})$ exist for all $\mathbf{x} \in A$ a Borel set contained in G . Also, $\mathcal{X}_A D\mathbf{h}(\mathbf{x})$, considered as a matrix has all of its entries Borel measurable if \mathbf{h} is continuous because these are obtained as partial derivatives which are limits of continuous functions. Of course this is automatic if \mathbf{h} is Lipschitz because then A is all of G other than a set of measure zero. Recall that if \mathbf{h} is Lipschitz continuous, $\mathbf{h}(E)$ is \mathcal{H}^n Hausdorff measurable whenever E is n dimensional Lebesgue measurable.

For convenience $0 < \varepsilon < 1/4$ in what follows.

For $\mathbf{x} \in A$, let $D\mathbf{h}(\mathbf{x}) \equiv R(\mathbf{x})U(\mathbf{x})$ where $R(\mathbf{x})$ preserves lengths and

$$U(\mathbf{x}) \equiv (D\mathbf{h}(\mathbf{x})^* D\mathbf{h}(\mathbf{x}))^{1/2}.$$

This is the right polar factorization of Theorem 1.5.5. Let A^+ denote those points of A for which $U(\mathbf{x})^{-1}$ exists, where $\det(D\mathbf{h}(\mathbf{x})^* D\mathbf{h}(\mathbf{x})) > 0$. Thus this is a Borel measurable subset of A . Note that $\det(D\mathbf{h}(\mathbf{x})^* D\mathbf{h}(\mathbf{x})) \geq 0$ for all $\mathbf{x} \in A$.

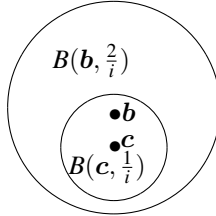
Let B be a Borel measurable subset of A^+ and let $\mathbf{b} \in B$. Let \mathcal{S} be a countable dense subset of the space of symmetric invertible matrices and let \mathcal{C} be a countable dense subset of B . The idea is to decompose B into countably many Borel sets E on which \mathbf{h} is one to one and Lipschitz with \mathbf{h}^{-1} Lipschitz on $\mathbf{h}(E)$. This will be done by establishing 17.6 given below where T is an invertible symmetric transformation. For $T \in \mathcal{S}$ and $\mathbf{c} \in \mathcal{C}$, define $E(T, \mathbf{c}, i)$ to be those $\mathbf{b} \in B(\mathbf{c}, \frac{1}{i})$ such that for all $\mathbf{a} \in B(\mathbf{b}, \frac{2}{i})$,

$$|\mathbf{h}(\mathbf{a}) - \mathbf{h}(\mathbf{b}) - D\mathbf{h}(\mathbf{b})(\mathbf{a} - \mathbf{b})| < \varepsilon |T(\mathbf{a} - \mathbf{b})| \quad (17.2)$$

and also $U(\mathbf{b})$ is close enough to T that the following hold.

$$\inf_{v \neq 0} \frac{|D\mathbf{h}(\mathbf{b})v|}{|Tv|} = \inf_{v \neq 0} \frac{|U(\mathbf{b})v|}{|Tv|} > 1 - \varepsilon, \quad (17.3)$$

$$\sup_{v \neq 0} \frac{|D\mathbf{h}(\mathbf{b})v|}{|Tv|} = \sup_{v \neq 0} \frac{|U(\mathbf{b})v|}{|Tv|} < 1 + \varepsilon \quad (17.4)$$



Note that it is not clear whether $c \in E(T, c, i)$ because of the above two requirements. What is going on here is that we are looking for b such that $Dh(b)$ is sufficiently close to one of those T which also are in a piece of B . Thus we start with one of those T and one of those points c and look for all b , if any, which do the right things. In one dimension, the T and $Dh(b)$ would just be slopes. There are countably many of these pieces of B being denoted as $E(T, c, i)$.

The union of these $E(T, c, i)$ is all of B because if $b \in B$,

$$|h(a) - h(b) - Dh(b)(a - b)| < \varepsilon |U(b)(a - b)| \quad (17.5)$$

whenever $a \in B(b, \frac{2}{i})$ provided i is sufficiently large. Thus also, by Lemma 17.2.1, there is $T \in \mathcal{S}$ such that the above holds for $U(b)$ replaced with T and $a \in B(b, \frac{2}{i})$ and also 17.3, 17.4. Thus $b \in E(T, c, i)$, so indeed the union of these sets is B .

Now let $a, b \in E(T, c, i)$. Since $a, b \in E(T, c, i)$, a, b are within $1/i$ of c and so a is within $2/i$ of b and so 17.2 holds because of the definition of $E(T, c, i)$. Therefore, from 17.2 and the inequalities which follow, 17.3 and 17.4,

$$(1 - 3\varepsilon) |T(a - b)| \leq |h(a) - h(b)| \leq (1 + 3\varepsilon) |T(a - b)| \quad (17.6)$$

Indeed, from 17.5, 17.4,

$$\begin{aligned} |h(a) - h(b)| &< |U(b)(a - b)| + \varepsilon |U(b)(a - b)| = (1 + \varepsilon) |U(b)(a - b)| \\ &< (1 + \varepsilon)^2 |T(a - b)| < (1 + 3\varepsilon) |T(a - b)| \end{aligned}$$

The other side of 17.6 is similar.

There are countably many of these $E(T, c, i)$ each being a Borel set. Therefore, B is a disjoint union of these sets called $\{E_k\}$ where I will denote the special T as T_k corresponding to E_k . Thus from 17.6 and the definition of Hausdorff measure, it follows that

$$\begin{aligned} \mathcal{H}^n(h(E_k)) &\in [(1 - 3\varepsilon) \mathcal{H}^n(T_k E_k), (1 + 3\varepsilon) \mathcal{H}^n(T_k E_k)] \\ &= [(1 - 3\varepsilon) m_n(T_k E_k), (1 + 3\varepsilon) m_n(T_k E_k)] \end{aligned} \quad (17.7)$$

From 17.3 and 17.4 and $b \in E_k$,

$$U(b)(B(\mathbf{0}, 1)) \subseteq (1 + \varepsilon) T_k(B(\mathbf{0}, 1)), U(b)(B(\mathbf{0}, 1)) \supseteq (1 - \varepsilon) T_k(B(\mathbf{0}, 1))$$

which implies on taking the Lebesgue measure that

$$(1 - \varepsilon)^n |\det(T_k)| \leq \det(U(b)) \leq (1 + \varepsilon)^n |\det(T_k)|$$

Therefore, from 17.7,

$$\begin{aligned} \mathcal{H}^n(h(E_k)) &\in [(1 - 3\varepsilon) |\det(T_k)| m_n(E_k), (1 + 3\varepsilon) |\det(T_k)| m_n(E_k)] \\ &= \left[\int_{E_k} (1 - 3\varepsilon) |\det(T_k)| dm_n, \int_{E_k} (1 + 3\varepsilon) |\det(T_k)| dm_n \right] \\ &\subseteq \left[\int_{E_k} (1 - 3\varepsilon) \left| \frac{\det(U(x))}{(1 + \varepsilon)^n} \right| dm_n(x), \int_{E_k} (1 + 3\varepsilon) \left| \frac{\det(U(x))}{(1 - \varepsilon)^n} \right| dm_n(x) \right] \end{aligned} \quad (17.8)$$

Note that this does not assume h is one to one on B .

Lemma 17.3.1 *Let $\mathbf{h} : G \rightarrow \mathbb{R}^m$ be Lipschitz. Let $B \subseteq A^+$ where B is Borel and where A consists of the points $\mathbf{x} \in G$ where $D\mathbf{h}(\mathbf{x})$ exists and A^+ consists of those points $\mathbf{x} \in A$ where for $D\mathbf{h}(\mathbf{x}) = R(\mathbf{x})U(\mathbf{x})$ in which $R^*R = I$, $\det(U(\mathbf{x})) > 0$. Then if \mathbf{h} is one to one on B ,*

$$\mathcal{H}^n(\mathbf{h}(B)) = \int_B \det(U(\mathbf{x})) dm_n(x) \quad (17.9)$$

Also for $Z \equiv A \setminus A^+$,

$$\mathcal{H}^n(\mathbf{h}(Z)) = 0 \quad (17.10)$$

regardless of whether \mathbf{h} is one to one. Letting $\#(\mathbf{y})$ be the number of points in $\mathbf{h}^{-1}(\mathbf{y}) \equiv \{\mathbf{x} \in G : \mathbf{h}(\mathbf{x}) = \mathbf{y}\} \cap B$, in the general case where \mathbf{h} is not required to be one to one,

$$\int_{\mathbf{h}(B)} \#(\mathbf{y}) d\mathcal{H}^n(y) = \int_B \det(U(\mathbf{x})) dm_n(x)$$

Proof: Let the $\{E_k, T_k\}$ be as described above where T_k goes with E_k . Let the union of these E_k be A^+ . Since the E_k are disjoint and \mathbf{h} is one to one, 17.9 follows from 17.8 applied to $B \cap E_k$ and summing over k , since ε is arbitrary.

Consider now the next assertion which is a form of Sard's lemma. Let P be the projection onto \mathbb{R}^m . Now consider 17.8 where we assume \mathbf{h} is Lipschitz on G . Since this holds for any small positive ε , it follows that

$$\mathcal{H}^n(\mathbf{h}(E_k \cap B)) = \int_{\mathbf{h}(B)} \mathcal{H}_{\mathbf{h}(E_k \cap B)}(\mathbf{y}) d\mathcal{H}^n(y) = \int_{E_k \cap B} |\det(U(\mathbf{x}))| dm_n(x) \quad (17.11)$$

First suppose G is bounded. Let $\mathbf{k}_\varepsilon(\mathbf{x}) \equiv \begin{pmatrix} \mathbf{h}(\mathbf{x}) \\ \varepsilon \mathbf{x} \end{pmatrix}$ so \mathbf{k}_ε is one to one. Then for all $\mathbf{x} \in A$, $\det(D\mathbf{k}_\varepsilon(\mathbf{x})^* D\mathbf{k}_\varepsilon(\mathbf{x})) = \det(D\mathbf{h}(\mathbf{x})^* D\mathbf{h}(\mathbf{x}) + \varepsilon^2 I_n) > 0$. Note that $P\mathbf{k}_\varepsilon = \mathbf{h}$ and for \mathbf{k}_ε , $A = A^+$. Letting $\{E_k\}$ be the disjoint Borel sets on which \mathbf{k}_ε is Lipschitz and one to one with inverse also Lipschitz, it follows

$$\mathcal{H}^n(\mathbf{k}_\varepsilon(Z \cap E_k)) = \int_{Z \cap E_k} \det(D\mathbf{h}(\mathbf{x})^* D\mathbf{h}(\mathbf{x}) + \varepsilon^2 I_n)^{1/2} dx$$

Also, since $\mathbf{h} = P\mathbf{k}_\varepsilon$, where P is Lipschitz with Lipschitz constant no more than 1, it follows from Lemma 17.1.2 that

$$\mathcal{H}^n(\mathbf{h}(Z \cap E_k)) \leq \int_{Z \cap E_k} \det(D\mathbf{h}(\mathbf{x})^* D\mathbf{h}(\mathbf{x}) + \varepsilon^2 I_n)^{1/2} dx$$

Then, $\mathbf{h}(Z) \subseteq \cup_k \mathbf{h}(Z \cap E_k)$. Hence,

$$\begin{aligned} \mathcal{H}^n(\mathbf{h}(Z)) &\leq \sum_k \mathcal{H}^n(\mathbf{h}(Z \cap E_k)) \leq \sum_k \int_{Z \cap E_k} \det(D\mathbf{h}(\mathbf{x})^* D\mathbf{h}(\mathbf{x}) + \varepsilon^2 I_n)^{1/2} dx \\ &= \int_Z \det(D\mathbf{h}(\mathbf{x})^* D\mathbf{h}(\mathbf{x}) + \varepsilon^2 I_n)^{1/2} dx \end{aligned}$$

Since \mathbf{h} is assumed Lipschitz, the expression in the integrand is bounded independent of ε and so, since G is bounded, the dominated convergence theorem applies and it follows that

$$\mathcal{H}^n(\mathbf{h}(Z)) \leq \int_Z \det(D\mathbf{h}(\mathbf{x})^* D\mathbf{h}(\mathbf{x}))^{1/2} dx = 0 \quad (17.12)$$

In case G is not bounded, apply the above to $G_n \equiv G \cap B(\mathbf{0}, n)$ for $n \in \mathbb{N}$. Then pass to a limit.

Now consider \mathbf{h} is only Lipschitz, maybe not one to one. Adding over k in 17.11,

$$\int_{\mathbf{h}(B)} \sum_k \mathcal{X}_{\mathbf{h}(E_k \cap B)}(\mathbf{y}) d\mathcal{H}^n(y) = \int_{\mathbf{h}(B)} \#(\mathbf{y}) d\mathcal{H}^n(y) = \int_B |\det(U(\mathbf{x}))| dm_n(x)$$

This is because if $\mathbf{x} \in B$ and $\mathbf{h}(\mathbf{x}) = \mathbf{y}$ then $\mathbf{x} \in B \cap E_k$ for some values of k but \mathbf{h} is one to one on $B \cap E_k$ and so this happens for at most one $\mathbf{x} \in E_k \cap B$.

Now suppose F is a Borel set in $\mathbf{h}(G)$ so $\mathbf{h}^{-1}(F)$ is a Borel set in \mathbb{R}^n . In the above let $B = \mathbf{h}^{-1}(F) \cap A^+$. Then

$$\int_{\mathbf{h}(\mathbf{h}^{-1}(F) \cap A^+)} \#(\mathbf{y}) d\mathcal{H}^n(y) = \int_{\mathbf{h}^{-1}(F) \cap A^+} |\det(U(\mathbf{x}))| dm_n(x)$$

Then this is

$$\begin{aligned} \int_{\mathbf{h}(A)} \mathcal{X}_F(\mathbf{y}) \#(\mathbf{y}) d\mathcal{H}^n(y) &= \int_{\mathbf{h}(A^+)} \mathcal{X}_F(\mathbf{y}) \#(\mathbf{y}) d\mathcal{H}^n(y) \\ &= \int_{A^+} \mathcal{X}_F(\mathbf{h}(\mathbf{x})) |\det(U(\mathbf{x}))| dm_n(x) \\ &= \int_A \mathcal{X}_F(\mathbf{h}(\mathbf{x})) |\det(U(\mathbf{x}))| dm_n(x) \end{aligned}$$

because $\mathbf{h}(A \setminus A^+)$ has \mathcal{H}^n measure zero and on $A \setminus A^+$, $|\det(U(\mathbf{x}))| = 0$. Since \mathbf{h} is Lipschitz, Rademacher's theorem implies that $G \setminus A$ has m_n measure zero and so also $\mathbf{h}(G \setminus A)$ has \mathcal{H}^n measure zero. Thus for F a Borel set,

$$\int_{\mathbf{h}(G)} \mathcal{X}_F(\mathbf{y}) \#(\mathbf{y}) d\mathcal{H}^n(y) = \int_G \mathcal{X}_F(\mathbf{h}(\mathbf{x})) |\det(U(\mathbf{x}))| dm_n(x)$$

For $\{E_k\}$ disjoint bounded Borel sets whose union is A^+ which are described above, consider $\lambda(E) \equiv \mathcal{H}^n(E \cap \mathbf{h}(E_k))$ for E an \mathcal{H}^n measurable set. This makes sense and is a measure because \mathbf{h} is one to one on E_k , $\mathbf{h}(E_k)$ is \mathcal{H}^n measurable because \mathbf{h} is Lipschitz on E_k and E_k is a Borel set, hence by Lemma 17.1.4 $\mathbf{h}(E_k)$ is \mathcal{H}^n measurable. λ is a finite measure because these E_k are all bounded and from what was shown above,

$$\lambda(\mathbb{R}^m) = \mathcal{H}^n(\mathbf{h}(E_k)) = \int_{E_k} \det(U(\mathbf{x})) dx < \infty$$

By Lemma 9.8.4, λ is regular on Borel sets. However, by Theorem 16.6.1, whenever E is a \mathcal{H}^n measurable set, there exists a Borel set F such that $\lambda(E) = \lambda(F)$ and $F \supseteq E$. Therefore, by Lemma 9.8.4, λ is a regular measure and if E is \mathcal{H}^n measurable, there exist F, H with these being Borel sets and such that $F \subseteq E \subseteq H$ and $\lambda(H \setminus F) = 0$. Therefore,

$$\mathcal{X}_F(\mathbf{h}(\mathbf{x})) \det(U(\mathbf{x})) \leq \mathcal{X}_E(\mathbf{h}(\mathbf{x})) \det(U(\mathbf{x})) \leq \mathcal{X}_H(\mathbf{h}(\mathbf{x})) \det(U(\mathbf{x}))$$

and

$$\int_{E_k} (\mathcal{X}_H(\mathbf{h}(\mathbf{x})) \det(U(\mathbf{x})) - \mathcal{X}_F(\mathbf{h}(\mathbf{x})) \det(U(\mathbf{x}))) dx = 0$$

so, $\mathcal{X}_H(\mathbf{h}(\mathbf{x})) \det(U(\mathbf{x})) - \mathcal{X}_F(\mathbf{h}(\mathbf{x})) \det(U(\mathbf{x})) = 0$ a.e. By completeness of Lebesgue measure, $\mathbf{x} \rightarrow \mathcal{X}_E(\mathbf{h}(\mathbf{x})) \det(U(\mathbf{x})) \mathcal{X}_{E_k}(\mathbf{x})$ is Lebesgue measurable and

$$\begin{aligned} \int_{\mathbf{h}(E_k)} \mathcal{X}_E(\mathbf{y}) d\mathcal{H}^n &= \int_{\mathbf{h}(E_k)} \mathcal{X}_F(\mathbf{y}) d\mathcal{H}^n = \int_{E_k} \mathcal{X}_F(\mathbf{h}(\mathbf{x})) \det(U(\mathbf{x})) d\mathbf{x} \\ &= \int_{E_k} \mathcal{X}_E(\mathbf{h}(\mathbf{x})) \det(U(\mathbf{x})) d\mathbf{x} \end{aligned}$$

Using the above argument, we can add these over k and obtain

$$\begin{aligned} \int_{\mathbf{h}(G)} \#(\mathbf{y}) \mathcal{X}_E(\mathbf{y}) d\mathcal{H}^n &= \int_{\mathbf{h}(A)} \#(\mathbf{y}) \mathcal{X}_E(\mathbf{y}) d\mathcal{H}^n = \int_{\mathbf{h}(A^+)} \#(\mathbf{y}) \mathcal{X}_E(\mathbf{y}) d\mathcal{H}^n \\ &= \int_{A^+} \mathcal{X}_E(\mathbf{h}(\mathbf{x})) \det(U(\mathbf{x})) d\mathbf{x} \\ &= \int_A \mathcal{X}_E(\mathbf{h}(\mathbf{x})) \det(U(\mathbf{x})) d\mathbf{x} \\ &= \int_G \mathcal{X}_E(\mathbf{h}(\mathbf{x})) \det(U(\mathbf{x})) d\mathbf{x} \end{aligned} \quad (17.13)$$

because $G \setminus A$ has measure zero and so does $\mathbf{h}(G \setminus A)$ and $\det(U(\mathbf{x})) = 0$ on $A \setminus A^+$. Also, from 17.10, $\mathbf{h}(A \setminus A^+)$ has measure zero. ■

Note that $\mathcal{H}^n(G \setminus A^+) = 0$ so it suffices to let $\#(\mathbf{y})$ simply be the number of points in $\mathbf{h}^{-1}(\mathbf{y})$. This has almost shown the following theorem.

Definition 17.3.2 To save on notation, I will denote $\det(U(\mathbf{x}))$ as $J_*(\mathbf{x})$.

Theorem 17.3.3 Suppose $\mathbf{h} : G \rightarrow \mathbb{R}^m$ is Lipschitz, G some open set, and let A be the Borel set on which $D\mathbf{h}$ exists with $m_p(G \setminus A) = 0$. (Rademacher's theorem). Then if $g \geq 0$ and is \mathcal{H}^n measurable,

$$\int_{\mathbf{h}(G)} \#(\mathbf{y}) g(\mathbf{y}) d\mathcal{H}^n = \int_G g(\mathbf{h}(\mathbf{x})) J_*(\mathbf{x}) dm_n.$$

and everything makes sense where here $\#(\mathbf{y})$ is defined as the number of times \mathbf{h} hits \mathbf{y} from points in A^+ or G . If \mathbf{h} is one to one on A^+ , we can replace $\#(\mathbf{y})$ with 1.

Proof: From 17.13 one can assert this holds for \mathcal{H}^n measurable simple functions and then, passing to a limit with monotone convergence theorem, one obtains the above theorem. ■

Next is an interesting version of the chain rule for Lipschitz maps. The proof of this theorem is based on the following lemma.

Lemma 17.3.4 If $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is Lipschitz, then if $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ for all $\mathbf{x} \in A$, then

$$\det(D\mathbf{h}(\mathbf{x})) = 0 \text{ a.e. } \mathbf{x} \in A$$

Proof: By the area formula, $0 = \int_{\{\mathbf{0}\}} \#(\mathbf{y}) d\mathbf{y} = \int_A |\det(D\mathbf{h}(\mathbf{x}))| d\mathbf{x}$, and so it follows that $\det(D\mathbf{h}(\mathbf{x})) = 0$ a.e. ■

Theorem 17.3.5 Let \mathbf{f}, \mathbf{g} be Lipschitz mappings from \mathbb{R}^n to \mathbb{R}^n with $\mathbf{g}(\mathbf{f}(\mathbf{x})) = \mathbf{x}$ on A , a measurable set. Then for a.e. $\mathbf{x} \in A$, $D\mathbf{g}(\mathbf{f}(\mathbf{x}))$, $D\mathbf{f}(\mathbf{x})$, and $D(\mathbf{g} \circ \mathbf{f})(\mathbf{x})$ all exist and $I = D(\mathbf{g} \circ \mathbf{f})(\mathbf{x}) = D\mathbf{g}(\mathbf{f}(\mathbf{x})) D\mathbf{f}(\mathbf{x})$.

Proof: By Lemma 17.3.4 there is a set of measure zero N_1 off which

$$\det(D(g \circ f)(x) - I) = 0$$

and in particular $D(g \circ f)(x)$ exists. Let N_2 be the set of measure zero off which f is differentiable. Let M be the set of points in $f(\mathbb{R}^n \setminus N_2)$ where, g fails to be differentiable. What about $f^{-1}(M)$? If $x \in f^{-1}(M)$ then $Dg(f(x))$ fails to exist and so x is in the first exceptional set N_1 or else in N_2 because $D(g \circ f)(x)$ will fail to exist. Thus $f^{-1}(M)$ is a set of measure zero. So let $x \notin N_1 \cup N_2$. Then for such x , $D(g \circ f)(x)$, $Dg(f(x))$, $Df(x)$ all exist and $I = Dg(f(x))Df(x)$. ■

You could give a generalization to the above by essentially repeating the argument.

Corollary 17.3.6 Suppose h is differentiable on A , a measurable set and that f, g are Lipschitz with $g(f(x)) = h(x)$ for $x \in A$. Then for a.e. $x \in A$,

$$Dh(x) = Dg(f(x))Df(x)$$

In other words, the chain rule holds off a set of measure zero.

17.4 The Divergence Theorem

Using Rademacher's theorem, all conditions are satisfied for Definition 14.3.1 provided each of the g_i used there are Lipschitz. When this happens, we say U is a bounded open set with a Lipschitz boundary which lies on one side of its boundary. Thus we obtain the divergence theorem. Here I will use \mathcal{H}^{p-1} for surface measure σ on the boundary of U since, by the area formula, this is what it is.

Theorem 17.4.1 Let U be a bounded open set with a Lipschitz boundary which lies on one side of its boundary. Then if $f \in C_c^1(\mathbb{R}^p)$,

$$\int_U f_{,k}(x) dm_p = \int_{\partial U} f n_k d\mathcal{H}^{p-1} \quad (17.14)$$

where $\mathbf{n} = (n_1, \dots, n_n)$ is the \mathcal{H}^{p-1} measurable unit outer normal. Also, if \mathbf{F} is a vector field such that each component is in $C_c^1(\mathbb{R}^p)$, then

$$\int_U \operatorname{div}(\mathbf{F}) dm_p = \int_{\partial U} \mathbf{F} \cdot \mathbf{n} d\mathcal{H}^{p-1}. \quad (17.15)$$

Proof: To obtain the first formula from the second which was proved earlier, consider

$$\mathbf{F} \equiv \begin{pmatrix} 0 & \cdots & f & \cdots & 0 \end{pmatrix}^T$$

where f is in the k^{th} slot. ■

You could approximate \mathbf{F} in the above theorem by convolving with a mollifier as in Lemma 16.3.1 yielding a modified \mathbf{F} with one in which each component is infinitely differentiable, apply the divergence theorem of Theorem 14.3.4 to these and pass to a limit using the dominated convergence theorem to obtain the divergence theorem for \mathbf{F} . Thus the following corollary is obtained.

Corollary 17.4.2 In the context of Theorem 17.4.1 it suffices to assume \mathbf{F} is Lipschitz.

17.5 The Coarea Formula

The area formula was discussed above. This formula implies that for E a measurable set

$$\mathcal{H}^n(\mathbf{f}(E)) = \int \mathcal{H}_E(\mathbf{x}) J_*(\mathbf{x}) dm$$

where $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ for \mathbf{f} a one to one Lipschitz mapping and $m \geq n$. The coarea formula is a statement about the Hausdorff measure of a set which involves the inverse image of \mathbf{f} . It looks a little like the method of shells in Calculus. We will let $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ where $m \leq n$ in what follows.

It is possible to obtain the coarea formula as a computation involving the area formula and some simple linear algebra and this is the approach taken here. I found this formula in [17] which has a somewhat different proof. I find this material very hard, so I hope what follows doesn't have grievous errors. I have never had occasion to use this coarea formula, but I think it is obviously of enormous significance and gives a very interesting geometric assertion. I will use the form of the chain rule in Theorem 17.3.5 as needed.

To begin with, here is the linear algebra identity. Recall that for a real matrix A^* is just the transpose of A . Thus AA^* and A^*A are symmetric.

Theorem 17.5.1 *Let A be an $m \times n$ matrix and let B be an $n \times m$ matrix for $m \leq n$. Then for I an appropriate size identity matrix, $\det(I + AB) = \det(I + BA)$.*

Proof: Use block multiplication to write

$$\begin{pmatrix} I+AB & 0 \\ B & I \end{pmatrix} \begin{pmatrix} I & A \\ 0 & I \end{pmatrix} = \begin{pmatrix} I+AB & A+ABA \\ B & BA+I \end{pmatrix}$$

$$\begin{pmatrix} I & A \\ 0 & I \end{pmatrix} \begin{pmatrix} I & 0 \\ B & I+BA \end{pmatrix} = \begin{pmatrix} I+AB & A+ABA \\ B & I+BA \end{pmatrix}$$

Hence

$$\begin{pmatrix} I+AB & 0 \\ B & I \end{pmatrix} \begin{pmatrix} I & A \\ 0 & I \end{pmatrix} = \begin{pmatrix} I & A \\ 0 & I \end{pmatrix} \begin{pmatrix} I & 0 \\ B & I+BA \end{pmatrix}$$

so

$$\begin{pmatrix} I & A \\ 0 & I \end{pmatrix}^{-1} \begin{pmatrix} I+AB & 0 \\ B & I \end{pmatrix} \begin{pmatrix} I & A \\ 0 & I \end{pmatrix} = \begin{pmatrix} I & 0 \\ B & I+BA \end{pmatrix}$$

which shows that the two matrices

$$\begin{pmatrix} I+AB & 0 \\ B & I \end{pmatrix}, \begin{pmatrix} I & 0 \\ B & I+BA \end{pmatrix}$$

are similar and so they have the same determinant. Thus $\det(I+AB) = \det(I+BA)$. Note that the two matrices are different sizes. ■

With these lemmas it is now possible to establish the coarea formula. First we define $\Lambda(n, m)$ as all possible ordered lists of m numbers taken from $\{1, 2, \dots, n\}$. Recall $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^m$ where $m \leq n$. Recall that this was part of the Binet Cauchy theorem, Theorem 1.9.14,

$$\det(D\mathbf{f}(\mathbf{x})D\mathbf{f}(\mathbf{x})^*) = \sum_{i \in \Lambda(n, m)} (\det D_{\mathbf{x}_i} \mathbf{f}(\mathbf{x}))^2$$

Now let $i_c \in \Lambda(n, n-m)$ consist of the remaining indices taken in order where $i \in \Lambda(n, m)$. For $i = (i_1, \dots, i_m)$, define $x_i \equiv (x_{i_1}, \dots, x_{i_m})$ and x_{i_c} to be the other components of x taken in order. Then let $f^i(x) \equiv \begin{pmatrix} f(x) \\ x_{i_c} \end{pmatrix}$. Thus there are $C(n, n-m) = C(n, m)$ different f^i featuring $C(n, m)$ different x_i, x_{i_c} .

Example 17.5.2 Say $f: \mathbb{R}^4 \rightarrow \mathbb{R}^2$. Here are some examples for f^i :

$$\begin{pmatrix} f_1(x_1, x_2, x_3, x_4) \\ f_2(x_1, x_2, x_3, x_4) \\ x_2 \\ x_4 \end{pmatrix}, \begin{pmatrix} f_1(x_1, x_2, x_3, x_4) \\ f_2(x_1, x_2, x_3, x_4) \\ x_1 \\ x_2 \end{pmatrix}, \begin{pmatrix} f_1(x_1, x_2, x_3, x_4) \\ f_2(x_1, x_2, x_3, x_4) \\ x_3 \\ x_4 \end{pmatrix}$$

Suppose first that $x_{i_c} = (x_{m+1} \dots x_n)^T$ so

$$f^i(x) = \begin{pmatrix} f(x) \\ x_{i_c} \end{pmatrix}, Df(x) = \begin{pmatrix} D_{x_i} f(x) & D_{x_{i_c}} f(x) \\ 0 & I \end{pmatrix}$$

and so from row operations, $\det Df^i(x) = \det D_{x_i} f(x)$. It is similar in the general case except one might have a sign change which is not important in what follows. So

$$\det Df^i(x) = \det D_{x_i} f(x). \quad (17.16)$$

Earlier with the area formula, we integrated $J_*(x) \equiv \det(Df(x)^* Df(x))^{1/2}$. With the coarea formula, we integrate $J^*(x) \equiv \det(Df(x) Df(x)^*)^{1/2}$. This proof involves doing this integration and seeing what happens. In case $n = m$ the claim of the theorem will follow from the area formula because $\mathcal{H}^0(E)$ is the number of elements of E , so one can assume if desired that in what follows $n > m$ although the argument does include this case.

Theorem 17.5.3 Let $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ where $n \geq m$ be a Lipschitz map. Let A be Lebesgue measurable. Then the following formula holds along with all measurability assertions needed for it to make sense.

$$\int_{\mathbb{R}^m} \mathcal{H}^{n-m}(A \cap f^{-1}(y)) dy = \int_{f(A)} \mathcal{H}^{n-m}(A \cap f^{-1}(y)) dy = \int_A J^*(x) dx \quad (17.17)$$

where $J^*(x) \equiv \det(Df(x) Df(x)^*)^{1/2}$.

Proof: First assume A is Borel, f differentiable on A . Now note that

$$\det(Df(x) Df(x)^*) = \sum_{i \in \Lambda(n, m)} \det(Df^i(x))^2$$

by the Binet Cauchy theorem and 17.16.

Lemma 17.5.4 Suppose A is a measurable nonempty set and $\det(Df(x) Df(x)^*) > 0$ for all $x \in A$. Then there exist disjoint, measurable A_i , one for each $i \in \Lambda(n, m)$ such that for all $j \neq i$, $\det(Df^j(x))^2 = 0$ for $x \in A_i$.

Proof: By assumption that $\det(Df(x)Df(x)^*) = \sum_{i \in \Lambda(n,m)} \det(Df^i(x))^2 > 0$, we can let $A_i \equiv \cap_{j \neq i} \{x : Df^j(x) = 0\}$. ■

Maybe some of these A_i are \emptyset but this will not matter.

Suppose f^i is one to one on a Borel set $E^i \subseteq \mathbb{R}^n$ which has positive measure and that its inverse, denoted as g^i is also Lipschitz on $f^i(E^i)$ and Df^i is invertible on E^i . Thus, for $x \in E^i \cap A$,

$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = y = f^i(x) = \begin{pmatrix} f(x) \\ x_{i_c} \end{pmatrix}, g_{i_c}^i(f^i(x)) = y_2 = x_{i_c}, y_1 = f(g^i(y)) \quad (17.18)$$

Differentiate $y_1 = f(g^i(y))$ with respect to y_2 to obtain

$$\begin{aligned} 0 &= D_{x_i} f(g^i(y)) D_{y_2} g_{i_c}^i(y) + D_{x_{i_c}} f(g^i(y)) D_{y_2} g_{i_c}^i(y) \\ &= D_{x_i} f(g^i(y)) D_{y_2} g_{i_c}^i(y) + D_{x_{i_c}} f(g^i(y)). \end{aligned} \quad (17.19)$$

Also,

$$Df^i(g^i(y)) Dg^i(y) = I, |\det(Dg^i(y))| = |\det Df^i(g^i(y))|^{-1}$$

Say $y = (y_1, y_2)^T$ and suppose $z = (z_i, z_{i_c})^T \in f^{-1}(y_1) \cap E^i \cap A$. Then

$$f^i \begin{pmatrix} z_i \\ z_{i_c} \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} f(z) \\ z_{i_c} \end{pmatrix}, \text{ so } \begin{pmatrix} z_i \\ z_{i_c} \end{pmatrix} = g^i \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$$

Thus,

$$g^i(f^{-1}(y_1) \cap E^i \cap A) = f^{-1}(y_1) \cap E^i \cap A \quad (17.20)$$

so if we fix y_1 , then $y_2 \rightarrow g^i(y_1, y_2)$ gives a parametrization for the Borel set $f^{-1}(y_1) \cap E^i \cap A$ and $D_{y_2} g_{i_c}^i(y) = I$.

Now from the area formula,

$$\int_{E^i \cap A} \det(Df^i(x)Df^i(x)^*)^{1/2} dx = \quad (17.21)$$

$$= \int_{f^i(E^i \cap A)} \det(Df^i(g^i(y))Df^i(g^i(y))^*)^{1/2} |\det Df^i(g^i(y))|^{-1} dy \quad (17.22)$$

Letting $y \equiv (y_1, y_2)$, and using what was just shown about $y_2 \rightarrow g^i(y_1, y_2)$ being a parametrization, the above integral can be expressed as the following iterated integral:

$$\int_{\mathbb{R}^m} \int_{f^{-1}(y_1) \cap E^i \cap A} \det(Df_{x_i}(g^i)Df_{x_i}(g^i)^*)^{1/2} |\det Df_{x_i}(g^i)|^{-1}(y_1, y_2) dy_2 dy_1, \quad (17.23)$$

Therefore, the inner integral is measurable and the integrand is

$$\det \left[\begin{pmatrix} D_{x_i} f(g^i(y)) & D_{x_{i_c}} f(g^i(y)) \end{pmatrix} \begin{pmatrix} D_{x_i} f(g^i(y))^* \\ D_{x_{i_c}} f(g^i(y))^* \end{pmatrix} \right]^{1/2} |\det Df_{x_i}(g^i(y))|^{-1}. \quad (17.24)$$

Let $A \equiv D_{x_i} f(g^i(y))$ so A is $m \times m$, $B \equiv D_{y_2} g_{i_c}^i(y)$ an $m \times (n-m)$. Using 17.19, 17.24 is of the form

$$\begin{aligned} &\det \left[\begin{pmatrix} A & -AB \end{pmatrix} \begin{pmatrix} A^* \\ -B^*A^* \end{pmatrix} \right]^{1/2} |\det A|^{-1} \\ &= \det[AA^* + ABB^*A^*]^{1/2} |\det A|^{-1} \\ &= \det[A(I + BB^*)A^*]^{1/2} |\det A|^{-1} = \det(I + BB^*)^{1/2} \end{aligned}$$

From Theorem 17.5.1, 17.23 equals $\det(I + B^*B)^{1/2}$. Note how the size of the matrices changes. Since $B = D_{y_2}g_i^i(y)$ and $D_{y_2}g_{i_c}^i(y) = I$, the above reduces to

$$\begin{aligned} \det(I + B^*B)^{1/2} &= \det \left[\begin{pmatrix} B^* & I \end{pmatrix} \begin{pmatrix} B \\ I \end{pmatrix} \right]^{1/2} = \\ \det \left[\begin{pmatrix} D_{y_2}g_i^i(y)^* & D_{y_2}g_{i_c}^i(y)^* \end{pmatrix} \begin{pmatrix} D_{y_2}g_i^i(y) \\ D_{y_2}g_{i_c}^i(y) \end{pmatrix} \right]^{1/2} &= \det(D_{y_2}g^i(y)^* D_{y_2}g^i(y))^{1/2} \end{aligned}$$

Therefore, from **area formula** and the above simplification of 17.24 and 17.20, 17.23

$$\begin{aligned} & \int \mathcal{X}_{E^i \cap A}(x) \det(D_{x_i}f(x) D_{x_i}f(x)^*)^{1/2} dx \\ &= \int_{\mathbb{R}^m} \int_{f^{-1}(y_1) \cap E^i \cap A} \det(D_{y_2}g^i(y)^* D_{y_2}g^i(y))^{1/2} dy_2 dy_1 \\ &= \int_{\mathbb{R}^m} \mathcal{H}^{n-m}(f^{-1}(y_1) \cap E^i \cap A) dy_1 \end{aligned} \quad (17.25)$$

Note how this also shows that $y_1 \rightarrow \mathcal{H}^{n-m}(f^{-1}(y_1) \cap E^i \cap A)$ is measurable since it equals the inner integral in an iterated integral having Borel integrand.

Now suppose that $A = A^+ \equiv \{x \in A : J^*(x) > 0\}$, and A is a Borel set. Let A_i be as in Lemma 17.5.4. Lemma 17.3.1 says there are disjoint Borel sets $\{E_j^i\}_{j=1}^\infty$ whose union is A_i on which the conditions for E^i in the above argument hold. Adding the above in 17.25 over j , for E^i replaced with E_j^i and using Lemma 17.5.4,

$$\begin{aligned} & \int \mathcal{X}_{A_i}(x) \det(Df^i(x) Df^i(x)^*)^{1/2} dx \\ &= \int \mathcal{X}_{A_i}(x) \det(Df(x) Df(x)^*)^{1/2} dx = \int_{\mathbb{R}^m} \mathcal{H}^{n-m}(f^{-1}(y_1) \cap A_i) dy_1 \end{aligned}$$

Now add the above over all A_i to obtain 17.17.

What if $A \setminus A^+ \neq \emptyset$? Then consider $(A \setminus A^+) \times \mathbb{R}^m \equiv \hat{A}$ as the new A and

$$\hat{f}(x_1, \dots, x_{n+m}) \equiv \begin{pmatrix} f(x_1, \dots, x_n) & \varepsilon x_{n+1} e_1 & \cdots & \varepsilon x_{n+m} e_m \end{pmatrix}$$

Thus $D\hat{f}(x) D\hat{f}(x)^* = Df(x) Df(x)^* + \varepsilon^{2m} I$ and so if $\varepsilon > 0$, $\det(D\hat{f}(x) D\hat{f}(x)^*) \equiv J_\varepsilon^*(x) \neq 0$ since $D\hat{f}(x) D\hat{f}(x)^*$ has positive eigenvalues at least ε^{2m} . Then from what was done above, letting \hat{E} be a bounded Borel set in \mathbb{R}^{n+m} and E the corresponding bounded set in \mathbb{R}^n ,

$$\begin{aligned} \int \mathcal{X}_{\hat{A} \cap \hat{E}}(x) J_\varepsilon^*(x) dm_{n+m} &= \int_{\mathbb{R}^m} \mathcal{H}^{n-m}(\hat{f}^{-1}(y_1) \cap \hat{A} \cap \hat{E}) dy_1 \\ &\geq \int_{\mathbb{R}^m} \mathcal{H}^{n-m}(f^{-1}(y_1) \cap (A \setminus A^+) \cap E) dy_1 \end{aligned}$$

where $\lim_{\varepsilon \rightarrow 0} J_\varepsilon^*(x) = 0$. This follows since projections decrease distance. Then by the dominated convergence theorem we can pass to a limit and find

$$\int_{\mathbb{R}^m} \mathcal{H}^{n-m}(f^{-1}(y_1) \cap (A \setminus A^+) \cap E) = 0.$$

Since this holds for any bounded E , this implies $\int_{\mathbb{R}^m} \mathcal{H}^{n-m}(\mathbf{f}^{-1}(\mathbf{y}_1) \cap (A \setminus A^+)) = 0$. Thus the desired formula holds in this case also.

Borel measurability of A can be replaced with Lebesgue measurability. Let A be Lebesgue measurable and let $F \subseteq A \subseteq G$ where $m_n(G \setminus F) = 0$ and F is F_σ and G is G_δ . From the above, $\int_{\mathbb{R}^m} \mathcal{H}^{n-m}(G \cap \mathbf{f}^{-1}(\mathbf{y})) dy = \int_{\mathbb{R}^m} \mathcal{H}^{n-m}(F \cap \mathbf{f}^{-1}(\mathbf{y})) dy$ and so from the above arguments, $\mathcal{H}^{n-m}(F \cap \mathbf{f}^{-1}(\mathbf{y})) = \mathcal{H}^{n-m}(G \cap \mathbf{f}^{-1}(\mathbf{y}))$ a.e. Thus $\mathbf{y} \rightarrow \mathcal{H}^{n-m}(A \cap \mathbf{f}^{-1}(\mathbf{y}))$ is measurable by completeness of the measure and

$$\begin{aligned} & \int_{\mathbb{R}^m} \mathcal{H}^{n-m}(A \cap \mathbf{f}^{-1}(\mathbf{y})) dy \\ & \in \left[\int_{\mathbb{R}^m} \mathcal{H}^{n-m}(F \cap \mathbf{f}^{-1}(\mathbf{y})) dy, \int_{\mathbb{R}^m} \mathcal{H}^{n-m}(G \cap \mathbf{f}^{-1}(\mathbf{y})) dy \right] \\ & = \left[\int_F J^*(\mathbf{x}) dx, \int_G J^*(\mathbf{x}) dx \right] = \left[\int_A J^*(\mathbf{x}) dx, \int_A J^*(\mathbf{x}) dx \right] \end{aligned}$$

We don't need to restrict A to be contained in N^C where N is the set of Lebesgue measure 0 where $D\mathbf{f}$ does not exist. Consider $N^C \cap B(0, k)$.

$$(N^C \cap B(0, k) \cap \mathbf{f}^{-1}(\mathbf{y})) \cup (N \cap B(0, k) \cap \mathbf{f}^{-1}(\mathbf{y})) = B(0, k) \cap \mathbf{f}^{-1}(\mathbf{y})$$

and the ends are \mathcal{H}^{n-m} measurable so it follows that so is $N \cap B(0, k) \cap \mathbf{f}^{-1}(\mathbf{y})$. Also $\mathbf{y} \rightarrow \mathcal{H}^{n-m}(N \cap B(0, k) \cap \mathbf{f}^{-1}(\mathbf{y}))$ is measurable by similar reasoning and

$$\begin{aligned} & \int_{\mathbb{R}^m} \mathcal{H}^{n-m}(N \cap B(0, k) \cap \mathbf{f}^{-1}(\mathbf{y})) dy \\ & = \int_{\mathbb{R}^m} \mathcal{H}^{n-m}(B(0, k) \cap \mathbf{f}^{-1}(\mathbf{y})) dy - \int_{\mathbb{R}^m} \mathcal{H}^{n-m}(N^C \cap B(0, k) \cap \mathbf{f}^{-1}(\mathbf{y})) dy = 0 \end{aligned}$$

So, passing to a limit and the monotone convergence theorem, we get

$$\int_{\mathbb{R}^m} \mathcal{H}^{n-m}(N \cap \mathbf{f}^{-1}(\mathbf{y})) dy = 0.$$

Therefore, the set of points where \mathbf{f} fails to be differentiable is irrelevant and can be ignored. ■

Also note that by definition,

$$\int_{\mathbb{R}^m} \mathcal{H}^{n-m}(A \cap \mathbf{f}^{-1}(\mathbf{y})) dy = \int_{\mathbf{f}(A)} \mathcal{H}^{n-m}(A \cap \mathbf{f}^{-1}(\mathbf{y})) dy.$$

Recall that $\mathcal{H}^0(E)$ equals the number of elements in E . Thus, if $n = m$, the coarea formula implies

$$\int_A J^* \mathbf{f}(\mathbf{x}) dx = \int_{\mathbf{f}(A)} \mathcal{H}^0(A \cap \mathbf{f}^{-1}(\mathbf{y})) dy = \int_{\mathbf{f}(A)} \#(\mathbf{y}) dy \geq \int_{\mathbf{f}(A)} 1 dy$$

Thus, this gives a version of Sard's theorem by letting the singular set S be A in the above.

17.6 Change of Variables

Now let $s(x) = \sum_{i=1}^p c_i \chi_{E_i}(x)$ where E_i is Lebesgue measurable and $c_i \geq 0$. Then

$$\begin{aligned} \int_{\mathbb{R}^n} s(x) J^*(f)(x) dx &= \sum_{i=1}^p c_i \int_{E_i} J^*(f)(x) dx = \sum_{i=1}^p c_i \int_{\mathbb{R}^m} \mathcal{H}^{n-m}(E_i \cap f^{-1}(y)) dy \\ &= \int_{\mathbb{R}^m} \sum_{i=1}^p c_i \mathcal{H}^{n-m}(E_i \cap f^{-1}(y)) dy = \int_{\mathbb{R}^m} \left[\int_{f^{-1}(y)} s d\mathcal{H}^{n-m} \right] dy \\ &= \int_{\mathbb{R}^m} \left[\int_{f^{-1}(y)} s d\mathcal{H}^{n-m} \right] dy = \int_{f(\mathbb{R}^n)} \left[\int_{f^{-1}(y)} s d\mathcal{H}^{n-m} \right] dy. \end{aligned} \quad (17.26)$$

Theorem 17.6.1 *Let $g \geq 0$ be Lebesgue measurable and let*

$$f : \mathbb{R}^n \rightarrow \mathbb{R}^m, \quad n \geq m, \quad f \text{ being Lipschitz}$$

Then

$$\int_{\mathbb{R}^n} g(x) J^*(f)(x) dx = \int_{f(\mathbb{R}^n)} \left[\int_{f^{-1}(y)} g(u) d\mathcal{H}^{n-m}(u) \right] dy. \quad (17.27)$$

Proof: Let $s_i \uparrow g$ where s_i is a simple function satisfying 17.26. Then let $i \rightarrow \infty$ and use the monotone convergence theorem to replace s_i with g . This proves the change of variables formula. ■

Note how if $m = n$ this will end up reducing to the conclusion of Theorem 11.10.2.

The following is an easy example of the use of the coarea formula to give a familiar relation.

Example 17.6.2 *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be given by $f(x) \equiv |x|$. Then $J^*(x)$ ends up being 1. Then by the coarea formula,*

$$\int_{B(\mathbf{0}, r)} dm_n = \int_0^r \mathcal{H}^{n-1}(B(\mathbf{0}, r) \cap f^{-1}(y)) dy = \int_0^r \mathcal{H}^{n-1}(\partial B(\mathbf{0}, y)) dy$$

Then $m_n(B(\mathbf{0}, r)) \equiv \alpha_n r^n = \int_0^r \mathcal{H}^{n-1}(\partial B(\mathbf{0}, y)) dy$. Then differentiate both sides to obtain $n\alpha_n r^{n-1} = \mathcal{H}^{n-1}(\partial B(\mathbf{0}, r))$. In particular $\mathcal{H}^2(\partial B(\mathbf{0}, r)) = 3\frac{4}{3}\pi r^2 = 4\pi r^2$. Of course α_n was computed earlier. Recall from Theorem 14.4.1 on Page 405

$$\alpha_n = \pi^{n/2} (\Gamma(n/2 + 1))^{-1}$$

Therefore, the $n - 1$ dimensional Hausdorff measure of the boundary of the ball of radius r in \mathbb{R}^n is $n\pi^{n/2} (\Gamma(n/2 + 1))^{-1} r^{n-1}$.

I think it is clear that you could generalize this to other more complicated situations. The above is nice because $J^*(x) = 1$. This won't be so in general when considering other level surfaces.

17.7 Integration and the Degree

There is a very interesting application of the degree to integration. I saw something like it in [20]. I want to consider the case where $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is only Lipschitz continuous, vanishing outside a bounded set. In the following proposition, let ϕ_ε be a symmetric nonnegative mollifier,

$$\phi_\varepsilon(\mathbf{x}) \equiv \frac{1}{\varepsilon^n} \phi\left(\frac{\mathbf{x}}{\varepsilon}\right), \text{ spt } \phi \subseteq B(\mathbf{0}, 1).$$

Ω will be a bounded open set. By Rademacher's theorem, \mathbf{h} satisfies $D\mathbf{h}(\mathbf{x})$ exists a.e. If U is a bounded open set, $\lim_{m \rightarrow \infty} D(\mathbf{h} * \psi_m) = \lim_{m \rightarrow \infty} D\mathbf{h} * \psi_m = D\mathbf{h}$ in $L^1(U; \mathbb{R}^{n \times n})$ where ψ_m is a mollifier. Thus a subsequence converges a.e.

Now recall the definition of the degree.

Definition 17.7.1 Let Ω be a bounded open set in \mathbb{R}^p and let $\mathbf{f} : \bar{\Omega} \rightarrow \mathbb{R}^p$ be continuous. Let $\mathbf{y} \notin \mathbf{f}(\partial\Omega)$. Then the degree is defined as follows: Let \mathbf{g} be infinitely differentiable,

$$\|\mathbf{f} - \mathbf{g}\|_{\infty, \bar{\Omega}} < \delta \equiv \text{dist}(\mathbf{f}(\partial\Omega), \mathbf{y}),$$

and \mathbf{y} is a regular value of \mathbf{g} . Then $\mathbf{y} \notin \mathbf{g}(\partial\Omega)$ and we define

$$d(\mathbf{f}, \Omega, \mathbf{y}) \equiv \sum \{ \text{sgn}(\det(D\mathbf{g}(\mathbf{x}))) : \mathbf{x} \in \mathbf{g}^{-1}(\mathbf{y}), \mathbf{x} \in \Omega \}$$

where the sum is finite by Lemma 15.1.5, defined to equal 0 if $\mathbf{g}^{-1}(\mathbf{y})$ is empty.

Also recall the fundamental integral identity.

Lemma 17.7.2 Let $\mathbf{y} \notin \mathbf{g}(\partial\Omega)$ for $\mathbf{g} \in C^\infty(\bar{\Omega}; \mathbb{R}^p)$. Also suppose \mathbf{y} is a regular value of \mathbf{g} . Then for all positive ε small enough,

$$\int_{\Omega} \phi_\varepsilon(\mathbf{g}(\mathbf{x}) - \mathbf{y}) \det D\mathbf{g}(\mathbf{x}) d\mathbf{x} = \sum \{ \text{sgn}(\det D\mathbf{g}(\mathbf{x})) : \mathbf{x} \in \mathbf{g}^{-1}(\mathbf{y}) \}$$

The sum is the definition of the degree for \mathbf{g} as described. There was also an important identity about homotopy Lemma 15.1.13 which includes the following.

Lemma 17.7.3 If \mathbf{h} is in $C^\infty(\bar{\Omega} \times [a, b], \mathbb{R}^p)$, and $\mathbf{0} \notin \mathbf{h}(\partial\Omega \times [a, b])$ then for $0 < \varepsilon < \text{dist}(\mathbf{0}, \mathbf{h}(\partial\Omega \times [a, b]))$, $t \rightarrow \int_{\Omega} \phi_\varepsilon(\mathbf{h}(\mathbf{x}, t)) \det D_1 \mathbf{h}(\mathbf{x}, t) d\mathbf{x}$ is constant for $t \in [a, b]$.

Let \mathbf{h}_m in what follows be $\mathbf{h} * \psi_m$ for \mathbf{h} continuous and ψ_m a mollifier. Eventually \mathbf{h} will be Lipschitz continuous. Let $\mathbf{y} \in \mathbf{h}(\partial\Omega)^C$ there is a ball $B(\mathbf{y}, \delta)$ such that $d(\mathbf{h}, \Omega, \mathbf{y}) = d(\mathbf{h}_m, \Omega, \hat{\mathbf{y}})$ for all m large enough and $\hat{\mathbf{y}} \in B(\mathbf{y}, \delta)$. This follows from the properties of the degree. Also, from a use of Lemma 17.7.2, this lemma gives

$$d(\mathbf{h}_m, \Omega, \mathbf{z}) = \int_{\Omega} \phi_\varepsilon(\mathbf{h}_m(\mathbf{x}) - \mathbf{z}) \det(D\mathbf{h}_m(\mathbf{x})) d\mathbf{x}$$

for all $\mathbf{z} \in B(\mathbf{y}, \delta)$ for sufficiently small ε . Now let $f \in C_c(B(\mathbf{y}, \delta))$. Then

$$\begin{aligned} \int f(\mathbf{z}) d(\mathbf{h}_m, \Omega, \mathbf{z}) d\mathbf{z} &= \int f(\mathbf{z}) \int_{\Omega} \phi_\varepsilon(\mathbf{h}_m(\mathbf{x}) - \mathbf{z}) \det(D\mathbf{h}_m(\mathbf{x})) d\mathbf{x} d\mathbf{z} \\ &= \int_{\Omega} \det(D\mathbf{h}_m(\mathbf{x})) \int f(\mathbf{z}) \phi_\varepsilon(\mathbf{h}_m(\mathbf{x}) - \mathbf{z}) d\mathbf{z} d\mathbf{x} \\ &= \int f(\mathbf{z}) \int_{\Omega} \det(D\mathbf{h}_m(\mathbf{x})) \phi_\varepsilon(\mathbf{h}_m(\mathbf{x}) - \mathbf{z}) d\mathbf{x} d\mathbf{z} \end{aligned}$$

So letting $\varepsilon \rightarrow 0$, $\int f(z) d(\mathbf{h}_m, \Omega, z) dz = \int f(z) \det(D\mathbf{h}_m(z)) dz$. Next suppose

$$f \in C_c(\mathbf{h}(\partial\Omega)^C)$$

Then $\text{spt}(f)$ is covered by finitely many of such balls like the above, $\{B_i\}_{i=1}^m$ with the property that if $g \in C_c(B_i)$, then

$$\int g(z) d(\mathbf{h}_m, \Omega, z) dz = \int_{\Omega} g(\mathbf{h}_m(x)) \det(D\mathbf{h}_m(x)) dx$$

Now let $\{\psi_j\}$ be a partition of unity on $\text{spt} f$ with $\text{spt} \psi_i \subseteq B_i$. Then

$$\begin{aligned} \int f(z) d(\mathbf{h}_m, \Omega, z) dz &= \int \sum_{i=1}^m \psi_i(z) f(z) d(\mathbf{h}_m, \Omega, z) dz \\ &= \sum_{i=1}^m \int \psi_i(z) f(z) d(\mathbf{h}_m, \Omega, z) dz = \sum_{i=1}^m \int_{\Omega} \psi_i(\mathbf{h}_m(x)) f(\mathbf{h}_m(x)) \det(D\mathbf{h}_m(x)) dx \\ &= \int_{\Omega} \sum_{i=1}^m \psi_i(\mathbf{h}_m(x)) f(\mathbf{h}_m(x)) \det(D\mathbf{h}_m(x)) dx = \int_{\Omega} f(\mathbf{h}_m(x)) \det(D\mathbf{h}_m(x)) dx \end{aligned}$$

If \mathbf{h} is Lipschitz, then $\lim_{m \rightarrow \infty} \mathbf{h}_m(x) = \mathbf{h}(x)$ uniformly and also for a.e. x , $D\mathbf{h}_m(x) \rightarrow D\mathbf{h}(x)$ and so, since everything is bounded, we can apply the dominated convergence theorem and conclude that

$$\int f(z) d(\mathbf{h}, \Omega, z) dz = \int_{\Omega} f(\mathbf{h}(x)) \det(D\mathbf{h}(x)) dx$$

This proves the following interesting proposition.

Proposition 17.7.4 *Let $f \in C_c(\mathbf{h}(\partial\Omega)^C)$ and let \mathbf{h} be Lipschitz on \mathbb{R}^n . Then*

$$\int f(\mathbf{y}) d(\mathbf{h}, \Omega, \mathbf{y}) dy = \int_{\Omega} f(\mathbf{h}(x)) \det(D\mathbf{h}(x)) dx.$$

Note that $d(\mathbf{h}, \Omega, \mathbf{y}) = 0$ if $\mathbf{y} \notin \mathbf{h}(\Omega)$ so the integral on the left is taken over $\mathbf{h}(\Omega)$. Recall that the area formula gives the formula

$$\int_{\mathbf{h}(\Omega)} f(\mathbf{y}) \#(\mathbf{y}) dy = \int_{\Omega} f(\mathbf{h}(x)) |\det(D\mathbf{h}(x))| dx.$$

You could probably say more. For example, the degree is constant on components of $\mathbf{h}(d\Omega)^C$ and so considering these components, you maybe could use the Riesz representation theorem for positive linear functionals to get this formula for more general f . Say Ω is open and connected, for example, and suppose $d(\mathbf{h}, \Omega, \mathbf{y})$ is 3 on $\mathbf{y} \in \mathbf{h}(\Omega)$. Then you could let $\Lambda f \equiv \int f(\mathbf{y}) 3 dy = \int_{\Omega} f(\mathbf{h}(x)) \det(D\mathbf{h}(x)) dx$ and this would be a positive linear functional on $C_c(\mathbf{h}(\partial\Omega)^C)$.

Chapter 18

Differential Forms

This is a generalization of single variable ideas from calculus to multiple dimensions. It generalizes the fundamental theorem of calculus of the form $\int_a^b f'(t) dt = f(b) - f(a)$ and integration by parts. The challenge is in keeping track of orientation. This was considered to some extent in the chapter on manifolds but differential forms allow for a more systematic presentation. Most of what follows involves whatever assumptions of smoothness are convenient, but to extend to the case of Lipschitz mappings, I will use Lemma 16.3.1.

Also note that the composition of Lipschitz mappings is Lipschitz and the usual chain rule will hold off some set of measure zero, see Theorem 17.3.5.

Recall from calculus that when you had a differential form, written as

$$\int_C a(x, y, z) dx + b(x, y, z) dy + c(x, y, z) dz$$

where C is an oriented curve, it gave the work done by the force field

$$\mathbf{F}(x, y, z) = (a(x, y, z), b(x, y, z), c(x, y, z))$$

on an object moving over the oriented curve C . How was it evaluated? You had a parametrization

$$\mathbf{r} : [a, b] \rightarrow \mathbb{R}^3, \mathbf{r}(t) = (x(t), y(t), z(t))$$

and then you did the following.

$$\int_a^b \left(a(x(t), y(t), z(t)) \frac{dx}{dt} + b(x(t), y(t), z(t)) \frac{dy}{dt} + c(x(t), y(t), z(t)) \frac{dz}{dt} \right) dt$$

Note how the orientation of the curve comes from the interval and the choice of parameterization. The interval $[a, b]$ is called the parameter domain and t is referred to as a parameter. The curve itself is some object in \mathbb{R}^3 .

You can think of the differential form $a(x, y, z) dx + b(x, y, z) dy + c(x, y, z) dz$ as a symbol which represents something which makes vector valued functions \mathbf{r} defined on $[a, b]$ into numbers according to the above procedure. You might also have written the line integral in the form

$$\int_{\mathbf{r}} a(x, y, z) dx + b(x, y, z) dy + c(x, y, z) dz$$

Thus it is a functional which makes vector valued functions defined on $[a, b]$ into numbers and if you change the orientation, this number changes. It is desired to extend this simple idea to functions of $k > 1$ variables. That which will take the place of an interval will be a box $\prod_{j=1}^k [a_j, b_j]$ or a chain of boxes formed by pasting boxes together along a common face. The latter will give a more general parameter domain than just a box. This is analogous to the notion of piecewise smooth curves where the curve was obtained by joining one smooth curve to another at an end point.

First here is some notation.

Notation 18.0.1 Let $\mathbf{r} : \prod_{j=1}^k [a_j, b_j] \rightarrow \mathbb{R}^p$, $p \geq k$ with $\mathbf{r} \in C^1 \left(\prod_{j=1}^k [a_j, b_j], \mathbb{R}^p \right)$. This last symbol means that \mathbf{r} is the restriction to $\prod_{j=1}^k [a_j, b_j]$ of a C^1 function defined on

all of \mathbb{R}^k . Now let I denote an ordered list of k indices taken from $\{1, 2, \dots, p\}$. Thus $I = (i_1, \dots, i_k)$. Then

$$\det \left(\frac{d\mathbf{r}^I}{d\mathbf{u}} \right) \equiv \det \begin{pmatrix} x_{i_1, u_1} & x_{i_1, u_2} & \cdots & x_{i_1, u_k} \\ x_{i_2, u_1} & x_{i_2, u_2} & \cdots & x_{i_2, u_k} \\ \vdots & \vdots & & \vdots \\ x_{i_k, u_1} & x_{i_k, u_2} & \cdots & x_{i_k, u_k} \end{pmatrix} \equiv \frac{\partial (x_{i_1}, \dots, x_{i_k})}{\partial (u_1, \dots, u_k)}$$

It is the same as $\det(D\mathbf{r}^I)$ where \mathbf{r}^I has values in \mathbb{R}^p and is obtained by keeping the rows of \mathbf{r} in the order determined by I and leaving out the other rows. More generally, suppose I is an ordered list of l indices and that J is an ordered list of l indices. Then

$$\det \left(\frac{d\mathbf{r}^I}{d\mathbf{u}_J} \right) \equiv \det \begin{pmatrix} x_{i_1, u_{j_1}} & x_{i_1, u_{j_2}} & \cdots & x_{i_1, u_{j_l}} \\ x_{i_2, u_{j_1}} & x_{i_2, u_{j_2}} & \cdots & x_{i_2, u_{j_l}} \\ \vdots & \vdots & & \vdots \\ x_{i_l, u_{j_1}} & x_{i_l, u_{j_2}} & \cdots & x_{i_l, u_{j_l}} \end{pmatrix} \equiv \frac{\partial (x_{i_1}, \dots, x_{i_l})}{\partial (u_{j_1}, \dots, u_{j_l})}, \quad \mathbf{x} = \mathbf{r}(\mathbf{u})$$

Now with this definition, here is the generalization of the differential forms defined in calculus.

Definition 18.0.2 A differential form of order k is $\omega \equiv \sum_I a_I(\mathbf{x}) d\mathbf{x}^I$ where

$$d\mathbf{x}^I \equiv dx_{i_1} \wedge dx_{i_2} \wedge \cdots \wedge dx_{i_k}$$

To save space, let $[\mathbf{a}, \mathbf{b}] \equiv \prod_{j=1}^k [a_j, b_j]$, $(\mathbf{a}, \mathbf{b}) \equiv \prod_{j=1}^k (a_j, b_j)$ etc. For $I = (i_1, \dots, i_k)$, $\int_{(\cdot)} \omega$ is a function mapping functions \mathbf{r} in $C^1(\prod_{j=1}^k [a_j, b_j], \mathbb{R}^p)$ or Lipschitz functions to \mathbb{R} , defined by

$$\int_{\mathbf{r}} \omega \equiv \int_{[\mathbf{a}, \mathbf{b}]} \sum_I a_I(\mathbf{r}(\mathbf{u})) \det \left(\frac{d\mathbf{r}^I(\mathbf{u})}{d\mathbf{u}} \right) d\mathbf{m}_k = \int_{[\mathbf{a}, \mathbf{b}]} \sum_I a_I(\mathbf{r}(\mathbf{u})) \det \left(\frac{d\mathbf{r}^I(\mathbf{u})}{d\mathbf{u}} \right) d\mathbf{m}_k$$

The sum is taken over all ordered lists of indices from $\{1, \dots, p\}$. Note that if there are any repeats in an ordered list I , then $\det \left(\frac{d\mathbf{r}^I(\mathbf{u})}{d\mathbf{u}} \right) = 0$ and so it suffices to consider the sum only over lists of indices in which there are no repeats. Thus the sum can be considered to consist of no more than $P(p, k)$ terms where this denotes the permutations of p things taken k at a time.

Consider the free Abelian group of mappings having some specified regularity from a given $[\mathbf{a}, \mathbf{b}]$ to \mathbb{R}^p . This consists of finite sums of the form $\sum m_{\mathbf{r}} \mathbf{r}$ and one can define

$$\int_{\sum m_{\mathbf{r}} \mathbf{r}} \omega \equiv \sum m_{\mathbf{r}} \int_{\mathbf{r}} \omega$$

Thus the integral defined on such an \mathbf{r} can be extended to give meaning to an arbitrary element of this free Abelian group.

Actually, if I, J are the same set of indices, listed in different order, then $\det \left(\frac{d\mathbf{r}^I(\mathbf{u})}{d\mathbf{u}} \right)$ will be $\pm \det \left(\frac{d\mathbf{r}^J(\mathbf{u})}{d\mathbf{u}} \right)$ so you could always write the differential form in terms of sums over

strictly increasing lists of indices, and if this is done, you can always have only $C(p, k) = \frac{p!}{k!(p-k)!}$ terms in the sum where this denotes combinations of p things taken k at a time.

It will always be assumed that a_I is as smooth as desired to make everything work. As to $\mathbf{r} \in C^1([a, b], \mathbb{R}^p)$ or Lipschitz, it is not required to have $D\mathbf{r}^I(\mathbf{u})$ be nonzero for any I . This means $\mathbf{r}([a, b])$ can be various sets which have points and edges.

Note that if $p < k$, then the functional $\sum_I a_I(\mathbf{x}) d\mathbf{x}^I$ should equal the zero function because you would have $x_{p+1} = \dots = x_k = 0$ and so there would be at least one row of zeros in the above determinant. Thus, I will suppose that $k \leq p$ in what follows.

Example 18.0.3 Consider the ball $B(\mathbf{0}, r)$. Spherical coordinates are $\mathbf{r} = (x, y, z)$ where

$$x = \rho \sin(\phi) \cos(\theta), y = \rho \sin(\phi) \sin(\theta), z = \rho \cos(\phi)$$

Let $(\rho, \phi, \theta) \in [0, r] \times [0, \pi] \times [0, 2\pi]$ which is a box like what was just described. This mapping is C^1 and onto the ball. However, it is clearly not one to one. However, \mathbf{r} is one to one on $(0, r] \times [0, \pi] \times [0, 2\pi]$ off a closed set S of m_3 measure zero. Also $\mathbf{r}(S)$ has measure zero by Sard's theorem. Recall from calculus that $\frac{\partial(x, y, z)}{\partial(\rho, \phi, \theta)} = \rho^2 \sin \phi$ which is positive if $\rho > 0$ and $\theta \in (0, \pi)$.

In a similar manner, you could obtain a solid ellipse.

$$x = a \sin(\phi) \cos(\theta), y = b \sin(\phi) \sin(\theta), z = c \cos(\phi)$$

for $a, b, c > 0$.

Example 18.0.4 Consider the boundary of the ball $B(\mathbf{0}, r)$. Spherical coordinates are

$$x = r \sin(\phi) \cos(\theta), y = r \sin(\phi) \sin(\theta), z = r \cos(\phi)$$

Let $(\phi, \theta) \in [0, \pi] \times [0, 2\pi]$ which is a box like what was just described. This mapping is C^1 and onto the boundary of this ball. It is clearly not one to one. However, off a set S of m_2 measure zero, the mapping is indeed one to one on what remains of this box and $\mathbf{r}(S)$ has \mathcal{H}^2 measure zero.

$$\frac{\partial(x, y)}{\partial(\phi, \theta)} = \begin{vmatrix} r \cos \phi \cos \theta & -r \sin \phi \sin \theta \\ r \cos \phi \sin \theta & r \sin \phi \cos \theta \end{vmatrix} = r^2 \cos \phi \sin \phi$$

$$\frac{\partial(x, z)}{\partial(\phi, \theta)} = \begin{vmatrix} r \cos \phi \cos \theta & -r \sin \phi \\ r \cos \phi \sin \theta & 0 \end{vmatrix} = r^2 \cos(\phi) \sin(\phi) \sin(\theta)$$

$$\frac{\partial(y, z)}{\partial(\phi, \theta)} = \begin{vmatrix} -r \sin \phi \sin \theta & -r \sin \phi \\ r \sin \phi \cos \theta & 0 \end{vmatrix} = r^2 \sin^2(\phi) \cos(\theta)$$

Similarly, you could obtain the boundary of an ellipse in the same way.

Example 18.0.5 One can obtain the cylinder of radius r of height h as follows.

$$x = r \cos(\theta), y = r \sin(\theta), z$$

where $(\theta, z) \in [0, 2\pi] \times [0, h]$ a box. The mapping is clearly C^1 but is not one to one.

More generally,

Example 18.0.6 Say $u \rightarrow (a(u), b(u))$ for $u \in [a, b]$ is a curve in the plane where a, b are C^1 and you consider

$$x = a(u), y = b(u), z = s$$

where $(s, u) \in [c, d] \times [a, b]$. This would be a surface in three dimensions. It may or may not be a one to one mapping depending on whether the curve is or is not a closed curve. It would not be one to one if the curve is a closed curve with $(x(b), y(b)) = (x(a), y(a))$.

Many other examples are available including higher dimensional versions of the above.

Lemma 18.0.7 Each differential form ω can be written in a unique way as

$$\sum_{i_1 < i_2 < \dots < i_k} a_{i_1, i_2, \dots, i_k}(\mathbf{x}) dx_{i_1} \wedge dx_{i_2} \wedge \dots \wedge dx_{i_k} \quad (18.1)$$

Also if σ is a permutation and $J = \sigma(I)$, then the following holds: $d\mathbf{x}^J = \text{sgn}(\sigma) d\mathbf{x}^I$

Proof: Consider the second claim. By definition,

$$\int_{\mathbf{r}} d\mathbf{x}^J \equiv \int_{[a, b]} \det\left(\frac{d\mathbf{x}^J}{d\mathbf{u}}\right) dm_k = \int_{[a, b]} \text{sgn}(\sigma) \det\left(\frac{d\mathbf{x}^I}{d\mathbf{u}}\right) dm_k \equiv \text{sgn}(\sigma) \int_{\mathbf{r}} d\mathbf{x}^I$$

this shows the second claim. The first claim that the form can be written as claimed is shown above. It involves switching rows and then combining terms.

Now it will be shown that the above sum in 18.1 is uniquely determined. This can be shown by considering a particular function in $C^1([a, b], \mathbb{R}^p)$ which is chosen auspiciously to reveal $a_{i_1, i_2, \dots, i_k}(\mathbf{x})$, $\mathbf{x} \in (a, b)$. Letting I be an ordered list of k indices from $\{1, \dots, p\}$. Say $I = (i_1, \dots, i_k)$ with $i_1 < i_2 < \dots < i_k$, an ascending list, meaning that the indices are increasing. Define $\mathbf{r}_\delta(\mathbf{u})$ as follows: $\mathbf{r}_\delta(\mathbf{u}) \equiv$

$$\begin{pmatrix} x_1 & \dots & x_{i_1} + \frac{\delta}{b_1 - a_1}(u_1 - a_1) & \dots & x_{i_k} + \frac{\delta}{b_k - a_k}(u_k - a_k) & \dots & x_p \end{pmatrix}^T \\ \equiv \mathbf{x} + \mathbf{d}_\delta(\mathbf{u})$$

So what is $\int_{\mathbf{r}_\delta} \omega$? By definition, it equals

$$\begin{aligned} \int_{\mathbf{r}_\delta} \omega &= \sum_{J \text{ ascending}} \int_{[a, b]} a_J(\mathbf{r}_\delta(\mathbf{u})) \det\left(\frac{d\mathbf{r}_{\delta J}(\mathbf{u})}{d\mathbf{u}}\right) dm_k \\ &= \prod_{j=1}^k \left(\frac{\delta}{b_j - a_j}\right) \int_{[a, b]} a_I(\mathbf{x} + \mathbf{d}_\delta(\mathbf{u})) dm_k(\mathbf{u}) \end{aligned} \quad (18.2)$$

because if J does not consist of the same indices as I , then there must be a row of zeros in $\frac{d\mathbf{r}_{\delta J}(\mathbf{u})}{d\mathbf{u}}$ and so the integral for that term in the above formula must equal 0. Re define u_j as $(u_j - a_j) \frac{\delta}{b_j - a_j} = u_j$, $j = 1, \dots, k$ so $\mathbf{u} \in [0, \delta] \equiv \prod_{s=1}^k [0, \delta]$. The appropriate Jacobian is then $\left(\prod_{j=1}^k \left(\frac{\delta}{b_j - a_j}\right)\right)^{-1}$ and so 18.2 reduces to $\int_{[0, \delta]} a_I(\mathbf{x} + \mathbf{u}) dm_k(\mathbf{u})$. Thus

$$\delta^{-k} \int_{\mathbf{r}_\delta} \omega = \delta^{-k} \int_{[0, \delta]} a_I(\mathbf{x} + \mathbf{u}) dm_k(\mathbf{u})$$

so by the fundamental theorem of calculus, $\lim_{\delta \rightarrow 0} \delta^{-k} \int_{\mathbf{r}_\delta} \omega = a_I(\mathbf{x})$. This verifies uniqueness. ■

18.1 The Wedge Product

Next is the definition of the wedge product.

Definition 18.1.1 Denote differential forms, the functionals defined above which act on $C^1([a, b], \mathbb{R}^p)$ as Ω^k . Then letting $\omega \in \Omega^k$ and $\eta \in \Omega^l$, let $\omega \wedge \eta \in \Omega^{k+l}$ be defined as follows. Letting

$$\omega \equiv \sum_{I \text{ ascending}} a_I(x) dx^I, \quad \eta \equiv \sum_{J \text{ ascending}} b_J(x) dx^J$$

then

$$\omega \wedge \eta \equiv \sum_{I \text{ ascending}} \sum_{J \text{ ascending}} a_I(x) b_J(x) dx^I \wedge dx^J$$

where for $I = (i_1, \dots, i_k), J = (j_1, \dots, j_l)$,

$$dx^I \wedge dx^J \equiv dx_{i_1} \wedge \dots \wedge dx_{i_k} \wedge dx_{j_1} \wedge \dots \wedge dx_{j_l}$$

This is well defined thanks to the above lemma which shows that there is only one way to write a differential form like the above sums in which I is ascending in each term.

What if I and J are not ascending? Does it still work? Let $I = (i_1, \dots, i_k), J = (j_1, \dots, j_l)$ and let $\sigma(I)$ be ascending and let $\eta(J)$ be ascending, σ, η being two permutations. Then from Lemma 18.0.7, the second claim,

$$\begin{aligned} dx^I \wedge dx^J & \stackrel{\substack{\in \Omega^k \\ \in \Omega^l}}{=} \stackrel{\substack{\text{writing each in ascending form} \\ \text{above definition}}}{=} \text{sgn}(\sigma) dx_{\sigma(I)} \wedge \text{sgn}(\eta) dx_{\eta(J)} \\ & \stackrel{\text{definition of the wedge product } dx_{\sigma(I)} \wedge dx_{\eta(J)}}{=} \text{sgn}(\sigma) \text{sgn}(\eta) dx_{\sigma(i_1)} \wedge \dots \wedge dx_{\sigma(i_k)} \wedge dx_{\eta(j_1)} \wedge \dots \wedge dx_{\eta(j_l)} \\ & = dx_{i_1} \wedge \dots \wedge dx_{i_k} \wedge dx_{j_1} \wedge \dots \wedge dx_{j_l} \end{aligned}$$

Thus it is correct to write $dx^I \wedge dx^J = dx_{i_1} \wedge \dots \wedge dx_{i_k} \wedge dx_{j_1} \wedge \dots \wedge dx_{j_l}$ even if $I = (i_1, \dots, i_k), J = (j_1, \dots, j_l)$ are not ascending.

This is the main idea behind the following fundamental result. In this lemma, I, J are not necessarily ascending.

Lemma 18.1.2 Let $\omega \equiv \sum_I a_I(x) dx^I, \eta \equiv \sum_J b_J(x) dx^J$ be in Ω^k and Ω^l respectively. Then

$$\omega \wedge \eta = \sum_{I, J} a_I(x) b_J(x) dx^I \wedge dx^J$$

Proof: For each I , let $\sigma_I(I)$ be ascending where σ_I is a permutation, similar for J . Then, as noted above in the proof of Lemma 18.0.7 and denoting by \hat{I}, \hat{J} the ascending lists of indices.

$$\omega = \sum_{\hat{I}} \left(\sum_{\{I: \sigma_I(I)=\hat{I}\}} a_I(x) \text{sgn}(\sigma_I) \right) dx^{\hat{I}}, \quad \eta = \sum_{\hat{J}} \left(\sum_{\{J: \sigma_J(J)=\hat{J}\}} b_J(x) \text{sgn}(\sigma_J) \right) dx^{\hat{J}}$$

Then by definition, $\omega \wedge \eta =$

$$\begin{aligned} & \sum_{\hat{f}, \hat{I}} \left(\sum_{\{I: \sigma_I(I)=\hat{I}\}} a_I(x) \operatorname{sgn}(\sigma_I) \right) \left(\sum_{\{J: \sigma_J(J)=\hat{f}\}} b_J(x) \operatorname{sgn}(\sigma_J) \right) dx^{\hat{I}} \wedge dx^{\hat{f}} \\ &= \sum_{\hat{f}, \hat{I}} \sum_{\{I: \sigma_I(I)=\hat{I}\}} \sum_{\{J: \sigma_J(J)=\hat{f}\}} a_I(x) b_J(x) \operatorname{sgn}(\sigma_I) (\operatorname{sgn}(\sigma_J)) dx^{\hat{I}} \wedge dx^{\hat{f}} \end{aligned}$$

By the alternating property of determinants,

$$= \sum_{\hat{f}, \hat{I}} \sum_{\{I: \sigma_I(I)=\hat{I}\}} \sum_{\{J: \sigma_J(J)=\hat{f}\}} a_I(x) b_J(x) dx^I \wedge dx^J = \sum_{I, J} a_I(x) b_J(x) dx^I \wedge dx^J \blacksquare$$

From this lemma, it is obvious that \wedge acts like multiplication in the sense that it is distributive over addition, and associative. However, it is not commutative. From properties of determinants,

$$\begin{aligned} & dx^I \wedge dx^J \\ &= dx_{i_1} \wedge \cdots \wedge dx_{i_k} \wedge dx_{j_1} \wedge \cdots \wedge dx_{j_l} \\ &= (-1)^{l+k-1} dx_{j_l} \wedge dx_{i_1} \wedge \cdots \wedge dx_{i_k} \wedge dx_{j_1} \wedge \cdots \wedge dx_{j_{l-1}} \\ &= (-1)^{l+k-1} (-1)^{l+k-1} dx_{j_{l-1}} \wedge dx_{j_l} \wedge dx_{i_1} \wedge \cdots \wedge dx_{i_k} \wedge dx_{j_1} \wedge \cdots \wedge dx_{j_{l-2}} \end{aligned}$$

etc. Thus you get eventually

$$\left((-1)^{l+k-1} \right)^l dx_{j_l} \wedge \cdots \wedge dx_{j_1} \wedge dx_{i_1} \wedge \cdots \wedge dx_{i_k}$$

Now $l^2 - l + lk = lk + l(l-1)$. However, $l(l-1)$ must be even. Therefore, this equals $(-1)^{|I||J|} dx^J \wedge dx^I$ where $|I|$ is the number of indices in I , similarly for J . This shows the following theorem which summarizes the algebraic properties of the wedge product.

Theorem 18.1.3 *Let α, β, γ be in some Ω^k . Then the following properties hold. If α, β are the same size, $(\alpha + \beta) \wedge \gamma = \alpha \wedge \gamma + \beta \wedge \gamma$. For α, β, γ arbitrary, then the following formula is obtained: $(\alpha \wedge \beta) \wedge \gamma = \alpha \wedge (\beta \wedge \gamma)$. Finally, there is the condition about what happens when order is reversed. If $\alpha \in \Omega^k$ and $\beta \in \Omega^l$, $\alpha \wedge \beta = (-1)^{lk} \beta \wedge \alpha$.*

Proof: The only claim which is not obvious is the last. However, this is also clear from Lemma 18.1.2 and the above computation involving $dx^I \wedge dx^J$. \blacksquare

18.2 The Exterior Derivative

A zero form is a function. Say f is such a smooth function. Then $df \in \Omega^1$ is defined as follows.

$$df \equiv \sum_{i=1}^p f_{x_i}(x) dx_i$$

Then for $\mathbf{r} \in C^1([0, 1], \mathbb{R}^p)$,

$$\int_{\mathbf{r}} df \equiv \int_0^1 \sum_{i=1}^p f_{x_i}(\mathbf{r}(u)) r'_i(u) du = \int_0^1 \frac{d}{du} (f(\mathbf{r}(u))) du = f(\mathbf{r}(1)) - f(\mathbf{r}(0))$$

In general, if you have something in Ω^k , you can define its exterior derivative as follows.

Definition 18.2.1 Let $\omega = \sum_{I \text{ ascending}} a_I(x) dx^I$. Then

$$d\omega \equiv \sum_{I \text{ ascending}} da_I(x) \wedge dx^I$$

It is clear that d is linear.

It doesn't matter whether ω is written in terms of ascending indices. The same formula holds with no change.

Lemma 18.2.2 Let $\omega = \sum_I a_I(x) dx^I$ then $d\omega = \sum_I da_I(x) \wedge dx^I$

Proof: Denote by \hat{I} the ascending indices. Then if $\sigma_I(I)$ is ascending,

$$\omega = \sum_{\hat{I}} \sum_{\{I\}=\{\hat{I}\}} a_I(x) dx^I = \sum_{\hat{I}} \left(\sum_{\{I\}=\{\hat{I}\}} a_I(x) \operatorname{sgn}(\sigma_I(I)) \right) dx^{\hat{I}}$$

Then since d is linear, $d\omega = \sum_{\hat{I}} \left(\sum_{\{I\}=\{\hat{I}\}} da_I(x) \operatorname{sgn}(\sigma_I(I)) \right) \wedge dx^{\hat{I}}$. Then since \wedge is also linear,

$$\begin{aligned} &= \sum_{\hat{I}} \left(\sum_{\{I\}=\{\hat{I}\}} da_I(x) \operatorname{sgn}(\sigma_I(I)) \wedge dx^{\hat{I}} \right) \\ &= \sum_{\hat{I}} \left(\sum_{\{I\}=\{\hat{I}\}} da_I(x) \wedge dx^I \right) = \sum_I da_I(x) \wedge dx^I \blacksquare \end{aligned}$$

Next is a product rule. First note that it follows right away from the definition and the product rule from beginning calculus that

$$d(fg) = d(f)g + gd(g)$$

Lemma 18.2.3 Let α, β be in Ω^k and Ω^l respectively. Then $d(\alpha \wedge \beta) = d\alpha \wedge \beta + (-1)^k \alpha \wedge d\beta$. Also $d^2 = 0$.

Proof: Let $\alpha = \sum_I a_I(x) dx^I, \beta = \sum_J b_J(x) dx^J$. Then $\alpha \wedge \beta = \sum_{I,J} a_I b_J dx^I \wedge dx^J$ and so $d(\alpha \wedge \beta)$ equals, thanks to the above lemma,

$$\begin{aligned} \sum_{I,J} d(a_I b_J) \wedge dx^I \wedge dx^J &= \sum_{I,J} [d(a_I(x)) b_J(x) + a_I(x) d(b_J(x))] \wedge dx^I \wedge dx^J \\ &= \sum_{I,J} d(a_I(x)) b_J(x) \wedge dx^I \wedge dx^J + \sum_{I,J} a_I(x) d(b_J(x)) \wedge dx^I \wedge dx^J \end{aligned}$$

From the definition of the wedge product and Lemma 18.1.2,

$$= d\alpha(x) \wedge \beta + \sum_{I,J} a_I(x) d(b_J(x)) \wedge dx^I \wedge dx^J$$

Now we will interchange the 1 form $d(b_J(x))$ and the k form dx^I . From Theorem 18.1.3

$$= d\alpha(x) \wedge \beta(x) + (-1)^k \sum_{I,J} a_I(x) dx^I \wedge d(b_J(x)) \wedge dx^J$$

Then from Lemma 18.1.2 again,

$$\begin{aligned} &= d\alpha(x) \wedge \beta(x) + (-1)^k \sum_I a_I(x) dx^I \wedge \sum_J d(b_J(x)) \wedge dx^J \\ &= d\alpha(x) \wedge \beta(x) + (-1)^k \alpha(x) \wedge d\beta(x) \end{aligned}$$

One of the important properties of the exterior derivative is that $d^2 = 0$. Let

$$\omega = \sum_I a_I(x) dx^I, \quad d\omega = \sum_I \sum_{r=1}^p a_{I,x_r} dx_r \wedge dx^I$$

Then by definition, $d^2\omega = \sum_I \sum_{r=1}^p \sum_{s=1}^p a_{I,x_r x_s} dx_s \wedge dx_r \wedge dx^I$

$$\begin{aligned} &= \sum_I \sum_{r < s} a_{I,x_r x_s} dx_s \wedge dx_r \wedge dx^I + \sum_I \sum_{s < r} a_{I,x_r x_s} dx_s \wedge dx_r \wedge dx^I \\ &= \sum_I \sum_{r < s} a_{I,x_r x_s} dx_s \wedge dx_r \wedge dx^I + \sum_I \sum_{r < s} a_{I,x_s x_r} dx_r \wedge dx_s \wedge dx^I \end{aligned}$$

In keeping with the convention that we assume a_I are as smooth as desired, we can conclude that the mixed partial derivatives are equal and so the above reduces to

$$\begin{aligned} &\sum_I \sum_{r < s} a_{I,x_r x_s} dx_s \wedge dx_r \wedge dx^I + \sum_I \sum_{r < s} a_{I,x_r x_s} dx_r \wedge dx_s \wedge dx^I \\ &= \sum_I \sum_{r < s} a_{I,x_r x_s} dx_s \wedge dx_r \wedge dx^I - \sum_I \sum_{r < s} a_{I,x_r x_s} dx_s \wedge dx_r \wedge dx^I = 0 \quad \blacksquare \end{aligned}$$

It might be interesting to note that if one means weak derivatives, then the mixed partial derivatives are always equal.

18.3 Stokes Theorem

Now that the algebra of differential forms has been presented, it is time for the main topic Stokes theorem. Recall $[a, b]$ is defined as $\prod_{l=1}^k [a_l, b_l]$. Let $x \in \mathbb{R}^p, p \geq k, r : [a, b] \rightarrow \mathbb{R}^p$, and let ω a $k-1$ form be given as follows

$$\omega = \sum_{I \in J} \alpha_I(x) dx_{i_1} \wedge \cdots \wedge dx_{i_{k-1}}$$

where here $I = (i_1, \dots, i_{k-1})$ is an increasing list of $k-1$ indices from $(1, 2, \dots, p)$ and J will denote the set of all such increasing lists of indices. Thus there are $C(p, k-1)$ elements in the set J . Assume that α_I is C^1 but here $x = r(u)$ is C^2 . After this case is done, it is easy to generalize.

This is a generalization of line integrals which is about integration over curves in space. Recall there was a parameter domain $[a, b]$ and a map $r : [a, b] \rightarrow \mathbb{R}^p$ and there were two orientations or directions over the curve which was the set of points $r([a, b])$. This concept of orientation is dealt with in multiple dimensions by the use of differential forms and the concept of determinants from linear algebra. The interval $[a, b]$ is replaced by $[a, b]$. Stokes theorem then is a statement about $r([a, b])$ and the boundary of $r([a, b])$ just as it is in the case of a line integral.

Stokes theorem relates the integral over some r to the integral over the “boundary” of r . This is defined next.

Definition 18.3.1 Denote as $\Lambda(k)$ the set of finite sums of differential forms of order k . I will now describe the boundary $\partial : \Lambda(k) \rightarrow \Lambda(k-1)$ by first defining it on \mathbf{r} and then, if desired, one would know it on all of $\Lambda(k)$. If $k = 1$, then $\partial \mathbf{r} \equiv \mathbf{r}(b) - \mathbf{r}(a)$. In general, for $\mathbf{r} : [\mathbf{a}, \mathbf{b}] \rightarrow \mathbb{R}^p$ and $l \leq k$,

$$\begin{aligned} & \partial_l \mathbf{r}(u_1, \dots, \hat{u}_l, \dots, u_k) \\ \equiv & \mathbf{r}(u_1, \dots, b_l, u_{l+1}, \dots, u_k) - \mathbf{r}(u_1, \dots, a_l, u_{l+1}, \dots, u_k) \equiv \mathbf{r}_{b_l} - \mathbf{r}_{a_l}. \end{aligned}$$

Note that $\mathbf{r}_{b_l}, \mathbf{r}_{a_l}$ are defined on

$$[a_1, b_1] \times \dots \times [a_{l-1}, b_{l-1}] \times [a_{l+1}, b_{l+1}] \times \dots \times [a_k, b_k] \equiv [\mathbf{a}, \mathbf{b}]_l.$$

Specifically, if $\omega = \sum_I a_I(\mathbf{x}) d\mathbf{x}^I$ where I denotes ordered lists of indices of length $k-1$ taken from $\{1, \dots, p\}$

$$\begin{aligned} \int_{\partial_l \mathbf{r}} \omega & \equiv \int_{[\mathbf{a}, \mathbf{b}]_l} \sum_I a_I(\mathbf{r}_{b_l}(\mathbf{u})) \det \left(\frac{d\mathbf{r}_{b_l}^I(\mathbf{u})}{d\mathbf{u}} \right) dm_{k-1} \\ & - \int_{[\mathbf{a}, \mathbf{b}]_l} \sum_I a_I(\mathbf{r}_{a_l}(\mathbf{u})) \det \left(\frac{d\mathbf{r}_{a_l}^I(\mathbf{u})}{d\mathbf{u}} \right) dm_{k-1} \end{aligned}$$

where \mathbf{u} comes from $[\mathbf{a}, \mathbf{b}]_l$. Here \mathbf{r}_{b_l} is \mathbf{r} with b_l in the l^{th} position, similar for \mathbf{r}_{a_l} . Thus $\det \left(\frac{d\mathbf{r}_{b_l}^I(\mathbf{u})}{d\mathbf{u}} \right), \det \left(\frac{d\mathbf{r}_{a_l}^I(\mathbf{u})}{d\mathbf{u}} \right)$ are $(-1)^{1+l} A_{1l}^J$ where A_{1l}^J is the $(1, l)^{\text{th}}$ cofactor of $\det(D\mathbf{r}^J)$ where J has length k and the matrix $D\mathbf{r}^J$ is the matrix of 18.4, having the top row

$$\begin{pmatrix} x_{j,u_1} & x_{j,u_2} & \dots & x_{j,u_k} \end{pmatrix}$$

Then $\int_{\partial \mathbf{r}} \omega \equiv \sum_I \int_{\partial_l \mathbf{r}} \omega$.

With this preparation, Stoke's theorem follows from a computation.

$$d\omega \equiv \sum_{I \in J} \sum_{j=1}^p \frac{\partial \alpha_I}{\partial x_j}(\mathbf{x}) dx_j \wedge dx_{i_1} \wedge \dots \wedge dx_{i_{k-1}}, \quad I = (i_1, \dots, i_{k-1})$$

Definition 18.3.2 As discussed earlier, define

$$\int_{\mathbf{r}} d\omega \equiv \sum_{I \in J} \sum_{j=1}^p \int_{[\mathbf{a}, \mathbf{b}]} \frac{\partial \alpha_I}{\partial x_j}(\mathbf{r}(\mathbf{u})) \frac{\partial (x_j, x_{i_1} \dots x_{i_{k-1}})}{\partial (u_1, \dots, u_k)} d\mathbf{u} \quad (18.3)$$

By definition, $\frac{\partial (x_j, x_{i_1} \dots x_{i_{k-1}})}{\partial (u_1, \dots, u_k)}$ is the determinant of

$$\begin{pmatrix} x_{j,u_1} & x_{j,u_2} & \dots & x_{j,u_k} \\ x_{i_1,u_1} & x_{i_1,u_2} & \dots & x_{i_1,u_k} \\ \vdots & \vdots & & \vdots \\ x_{i_{k-1},u_1} & x_{i_{k-1},u_2} & \dots & x_{i_{k-1},u_k} \end{pmatrix} \quad (18.4)$$

Note how this matrix is just the matrix of $D\mathbf{f}$ where

$$\mathbf{f}(\mathbf{u}) = (x_j(\mathbf{u}), x_{i_1}(\mathbf{u}), \dots, x_{i_{k-1}}(\mathbf{u}))^T.$$

Then expanding the determinant in 18.3 along the first row, it equals

$$= \sum_{I \in J} \sum_{j=1}^p \int_{[a,b]} \frac{\partial \alpha_I}{\partial x_j}(\mathbf{r}(\mathbf{u})) \overset{\text{expanding determinant}}{\sum_{l=1}^k \frac{\partial x_j}{\partial u_l} A_{1l}^I} d\mathbf{u}$$

where A_{1l}^I is the $(1, l)^{th}$ cofactor for the determinant of 18.4.

$$A_{1l}^I = (-1)^{1+l} \frac{\partial (x_{i_1}, \dots, x_{i_{k-1}})}{\partial (u_1, \dots, \hat{u}_l, \dots, u_k)}, I = (i_1, \dots, i_{k-1}) \quad (18.5)$$

Then this equals

$$= \sum_{I \in J} \sum_{l=1}^k \int_{[a,b]} \sum_{j=1}^p \frac{\partial \alpha_I}{\partial x_j}(\mathbf{r}(\mathbf{u})) \frac{\partial x_j}{\partial u_l} A_{1l}^I d\mathbf{u} = \sum_{I \in J} \sum_{l=1}^k \int_{[a,b]} \frac{\partial \alpha_I(\mathbf{r}(\mathbf{u}))}{\partial u_l} A_{1l}^I d\mathbf{u}$$

Now

$$\begin{aligned} \sum_I \frac{\partial \alpha_I(\mathbf{r}(\mathbf{u}))}{\partial u_l} A_{1l}^I &= \sum_I \frac{\partial}{\partial u_l} (\alpha_I(\mathbf{r}(\mathbf{u})) A_{1l}^I) - \sum_I \alpha_I(\mathbf{r}(\mathbf{u})) A_{1l,l}^I \\ &= \sum_I \frac{\partial}{\partial u_l} (\alpha_I(\mathbf{r}(\mathbf{u})) A_{1l}^I) \end{aligned}$$

By Lemma 7.11.2, that cofactor identity depending on equality of mixed partials. Therefore, from 18.3 and Fubini's theorem,

$$\begin{aligned} \int_{\mathbf{r}} d\omega &= \sum_{I \in J} \sum_{l=1}^k \int_{[a,b]} \frac{\partial}{\partial u_l} (\alpha_I(\mathbf{r}(\mathbf{u})) A_{1l}^I) \\ &= \sum_{I \in J} \sum_{l=1}^k \int_{[a,b]_l} \int_{[a_l, b_l]} \frac{\partial}{\partial u_l} ((\alpha_I(\mathbf{r}(\mathbf{u})) A_{1l}^I)) du_l du_I \\ &= \sum_{l=1}^k \int_{[a,b]_l} \sum_{I \in J} ((\alpha_I \circ \mathbf{r}) A_{1l}^I)(\mathbf{u}_I(b_l)) - ((\alpha_I \circ \mathbf{r}) A_{1l}^I)(\mathbf{u}_I(a_l)) du_I \end{aligned} \quad (18.6)$$

where here $[a, b]_l$ means the $[a_l, b_l]$ is missing in the product $[a, b]$ and $\mathbf{u}_I(b_l)$ is given by the formula $(u_1, \dots, u_{l-1}, b_l, u_{l+1}, \dots, u_k)$ with $\mathbf{u}_I(a_l)$ defined similarly. The term A_{1l}^I is the cofactor in 18.5 $(-1)^{1+l} \frac{\partial (x_{i_1}, \dots, x_{i_{k-1}})}{\partial (u_1, \dots, \hat{u}_l, \dots, u_k)}$. The term $\int_{[a,b]_l} ((\alpha_I \circ \mathbf{r}) A_{1l}^I)(\mathbf{u}_I(b_l)) du_I$ is an integration over the variables corresponding to a face of $[a, b]$ and so it is a kind of boundary term. By Definition 18.3.1 or simply making a definition that this is what we mean by the integral over the boundary, this is $\int_{\partial \mathbf{r}} \omega$. Thus, this proves Stokes' theorem.

Theorem 18.3.3 Let $\omega = \sum_I \alpha_I(\mathbf{x}) dx_{i_1} \wedge \dots \wedge dx_{i_{k-1}}$ be a $k-1$ form. Let $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^p, p \geq k$ be in $C^2([a, b]; \mathbb{R}^p)$. Then $\int_{\partial \mathbf{r}} \omega = \int_{\mathbf{r}} d\omega$.

Note that there is no assumption that $D\mathbf{r}$ has nonzero determinant. Everything above is valid under an assumption that \mathbf{r} is only C^2 . There was a reason why in calculus smooth curves had a parametrization with nonvanishing derivative. If the derivative vanishes, this

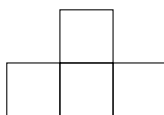
can yield a pointy place in the curve resulting from the given parametrization so it would not deserve to be called a smooth curve. It is the same here. Since $D\mathbf{r}$ is allowed to vanish, one can have $\mathbf{r}([a, b])$ many different kinds of sets. However, if you insist that $D\mathbf{r}$ be invertible, the points on the box would be preserved by an application of the implicit function theorem.

The above could be improved by using the above to approximate functions which are not C^2 with functions which are, obtained by mollifying, and then passing to a limit to get more general situations. In particular, consider \mathbf{r} a function in $C^1([a, b]; \mathbb{R}^p)$. Then by Lemma 16.3.1 there is a sequence of functions $\{\mathbf{r}_n\}$ each C^2 which converges uniformly to \mathbf{r} and such that $D\mathbf{r}_n$ converges uniformly to $D\mathbf{r}$ on $[a, b]$. Then Stoke's theorem holds with \mathbf{r} replaced with \mathbf{r}_n and so, passing to a limit as $n \rightarrow \infty$ one obtains Stoke's theorem for \mathbf{r} . This yields the following corollary.

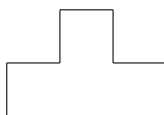
Corollary 18.3.4 *Let $\omega = \sum_I \alpha_I(\mathbf{x}) dx_{i_1} \wedge \cdots \wedge dx_{i_{k-1}}$ be an $k-1$ form, each α_I being $C^1([a, b])$. Let $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^p, p \geq k$ be in $C^1([a, b]; \mathbb{R}^p)$. Then $\int_{\partial \mathbf{r}} \omega = \int_{\mathbf{r}} d\omega$.*

Proof: From the above, $\int_{\partial \mathbf{r}_n} \omega = \int_{\mathbf{r}_n} d\omega$ where \mathbf{r}_n is C^2 . Both terms involve integrals over $[a, b]$ or $[a, b]_I$ and the convergence is uniform, so one can pass to a limit as $n \rightarrow \infty$ retaining the same formula with \mathbf{r}_n replaced with \mathbf{r} . ■

Suppose you have two boxes $[a, b]$ and $[c, d]$ and these intersect on a common face, say the l^{th} face. Thus $c_l = b_l$. Then in the above description for the boundary integral on the common face, the two contributions cancel because you have the same thing except one has $\left| \frac{b_l}{a_l} \right|$ and the other has $\left| \frac{d_l}{b_l} \right|$. Therefore, you would get Stokes theorem for the box consisting of these two pasted together, the boundary integrals consisting of the sum of the boundary integrals of the remaining faces. Continuing this way, consider a chain of these boxes such that each box intersects another along a complete face. Then you could do the above for each pair of boxes in the chain and note that the boundary integrals will cancel along the common faces. Thus, in place of a single box you could have a much more complicated shape and the boundary integral would take place over exactly those faces which do not have intersection with faces of other boxes whereas $\int_{\mathbf{r}} d\omega$ would take place over the union of the boxes since the boundaries are sets of measure zero relative to m_k . The following picture illustrates what is meant.



Thus, adding over the boxes yields a parameter domain which looks like the following for $\int_{\mathbf{r}} d\omega$.



By now, it should be clear that fairly general regions can be included. Also, we only need to have a_l continuous on the face of two of these intersecting boxes. This is an analog of a piece-wise smooth curve.

18.4 Lipschitz Maps

This will be based on the approximation with C^1 maps. Let $\mathbf{r} : [\mathbf{a}, \mathbf{b}] \subseteq \mathbb{R}^k \rightarrow \mathbb{R}^p$ where \mathbf{r} is Lipschitz and $k \leq p$ as before. The idea is to extend the above Stokes theorem to this case where \mathbf{r} is Lipschitz rather than C^1 . First extend \mathbf{r} as follows. For $t \geq b_l$ and $\mathbf{u}_l \in [\mathbf{a}, \mathbf{b}]_l$, let

$$\mathbf{r}(u_1, \dots, u_{l-1}, t, u_{l+1}, \dots, u_k) \equiv \mathbf{r}(u_1, \dots, u_{l-1}, b_l, u_{l+1}, \dots, u_k)$$

and also if $t \leq a_l$ and $\mathbf{u}_l \in [\mathbf{a}, \mathbf{b}]_l$, let

$$\mathbf{r}(u_1, \dots, u_{l-1}, t, u_{l+1}, \dots, u_k) \equiv \mathbf{r}(u_1, \dots, u_{l-1}, a_l, u_{l+1}, \dots, u_k)$$

This is done for each l . Then define for $h > 0$

$$\mathbf{r}^h(\mathbf{u}) \equiv \left(\frac{1}{2h}\right)^k \int_{-2h+u_1+\frac{2h}{(b_1-a_1)}(u_1-a_1)}^{u_1+\frac{2h}{(b_1-a_1)}(u_1-a_1)} \cdots \int_{-2h+u_k+\frac{2h}{(b_k-a_k)}(u_k-a_k)}^{u_k+\frac{2h}{(b_k-a_k)}(u_k-a_k)} \mathbf{r}(\mathbf{t}) dt_k \cdots dt_1 \quad (18.7)$$

Consider those integrals. When $u_1 = a_1$ you are integrating over $[a_1 - 2h, a_1]$ and when $u_1 = b_1$ you are integrating over $[b_1, b_1 + 2h]$. Of course it is similar for the other $[a_l, b_l]$. In general, each iterated integral is taken over an interval of length $2h$. For example,

$$u_1 + \frac{2h}{(b_1-a_1)}(u_1-a_1) - \left(-2h + u_1 + \frac{2h}{(b_1-a_1)}(u_1-a_1)\right) = 2h$$

Now consider the right end of the l^{th} interval in $[\mathbf{a}, \mathbf{b}]$ where $u_l = b_l$. Then this is describing one of the two faces corresponding to the l^{th} interval. By Fubini's theorem and the construction of the extension of \mathbf{r} , 18.7 implies that this equation simplifies to

$$\begin{aligned} \mathbf{r}^h(\mathbf{u}_l(b_l)) &= \left(\frac{1}{2h}\right)^{k-1} \int_{-2h+u_1+\frac{2h}{(b_1-a_1)}(u_1-a_1)}^{u_1+\frac{2h}{(b_1-a_1)}(u_1-a_1)} \cdots \int_{-2h+u_k+\frac{2h}{(b_k-a_k)}(u_k-a_k)}^{u_k+\frac{2h}{(b_k-a_k)}(u_k-a_k)} \\ &\quad \mathbf{r}(t_1, \dots, b_l, \dots, t_k) dt_k \cdots \widehat{dt_l} \cdots dt_1 \end{aligned}$$

Now by the fundamental theorem of calculus, for $i \neq l$,

$$\begin{aligned} \mathbf{r}_{u_i}^h(\mathbf{u}_l(b_l)) &= \left(\frac{1}{2h}\right)^{k-2} \int_{-2h+u_1+\frac{2h}{(b_1-a_1)}(u_1-a_1)}^{u_1+\frac{2h}{(b_1-a_1)}(u_1-a_1)} \cdots \int_{-2h+u_k+\frac{2h}{(b_k-a_k)}(u_k-a_k)}^{u_k+\frac{2h}{(b_k-a_k)}(u_k-a_k)} \\ &\quad \mathbf{r}\left(t_1, \dots, u_i + \frac{2h}{(b_i-a_i)}(u_i-a_i), \dots, b_l, \dots, t_k\right) \frac{1}{2h} \left(1 + \frac{2h}{b_i-a_i}\right) - \\ &\quad \mathbf{r}\left(t_1, \dots, -2h + u_i + \frac{2h}{(b_i-a_i)}(u_i-a_i), \dots, b_l, \dots, t_k\right) \frac{1}{2h} \left(1 + \frac{2h}{b_i-a_i}\right) \cdot \\ &\quad dt_k \cdots \widehat{dt_i} \cdots \widehat{dt_l} \cdots dt_1 \end{aligned}$$

By the material on Rademacher's theorem, that integrand is

$$\left(1 + \frac{2h}{b_i-a_i}\right) \frac{1}{2h} \int_{-2h+u_i+\frac{2h}{(b_i-a_i)}(u_i-a_i)}^{u_i+\frac{2h}{(b_i-a_i)}(u_i-a_i)} \mathbf{r}_{u_i}(t_1, \dots, t_i, \dots, b_l, \dots, t_k) dt_i$$

Now a.e. point of $[a, b]_l$ is a Lebesgue point and so we can pass to a limit and obtain pointwise a.e. convergence of $r_{u_l}^h(u_l(b_l))$ to $r_{u_l}(u_l(b_l))$.

Some comment on this might be useful because these Lebesgue points are not always at the center of the box of sides of length $2h$ determined by the limits of the iterated integrals. A given point $u_l(b_l) \in B(u_l(b_l), 3h)$ where this ball is taken with respect to $\|\cdot\|_\infty$. Thus, from Lebesgue's fundamental theorem of calculus, we have on this face at Lebesgue points the following converges to 0 as $h \rightarrow 0$ which is what was desired

$$\left(1 + \frac{2h}{b_l - a_l}\right) \left(\frac{3}{2}\right)^{k-1} \cdot \left(\frac{1}{3h}\right)^{k-1} \int_{B(u_l(b_l), 3h)} \|r_{u_l}(u_l(b_l)) - r_{u_l}(t_1, \dots, t_i, \dots, b_l, \dots, t_k)\| dt_l$$

Indeed, $\left(\frac{3}{2}\right)^{k-1} \left(\frac{1}{3h}\right)^{k-1} = \left(\frac{1}{2h}\right)^{k-1}$ and the integral in the above is taken over the larger set $B(u_l(b_l), 3h)$.

As to points of $[a, b]$, the pointwise convergence of $r_{u_j}^h(u)$ to $r_{u_j}(u)$ follows from similar reasoning but is a little less involved. Now r^h is clearly C^1 and so we have Stokes theorem for r^h .

$$\begin{aligned} & \sum_{l=1}^k \int_{[a, b]_l} \sum_{I \in J} \left((\alpha_I \circ r^h) A_{1l}^{lh} \right) (u_l(b_l)) - \left((\alpha_I \circ r) A_{1l}^{lh} \right) (u_l(a_l)) du_l \\ &= \int_r d\omega \equiv \sum_{I \in J} \sum_{j=1}^p \int_{[a, b]} \frac{\partial \alpha_I}{\partial x_j} (r^h(u)) \frac{\partial (x_j, x_{i_1} \cdots x_{i_{k-1}})}{\partial (u_1, \dots, u_k)} du \end{aligned}$$

where the superscript indicates that all is defined in terms of r^h . Then from the dominated convergence theorem, it follows that we can pass to a limit and obtain Stokes theorem where the boundary terms are defined from Rademacher's theorem on the faces of $[a, b]$.

Theorem 18.4.1 *Let $r : [a, b] \subseteq \mathbb{R}^k \rightarrow \mathbb{R}^p$, $p \geq k$ be Lipschitz and also suppose that $\alpha_I \in C^1(r([a, b]))$ for $I \subseteq J$, the set of increasing lists of $k-1$ indices from $(1, 2, \dots, p)$. Then one obtains Stokes theorem*

$$\begin{aligned} \int_{r(\partial[a, b])} \omega &\equiv \sum_{l=1}^k \int_{[a, b]_l} \sum_{I \in J} \left((\alpha_I \circ r) A_{1l}^I \right) (u_l(b_l)) - \left((\alpha_I \circ r) A_{1l}^I \right) (u_l(a_l)) du_l \\ &= \int_r d\omega \equiv \sum_{I \in J} \sum_{j=1}^p \int_{[a, b]} \frac{\partial \alpha_I}{\partial x_j} (r(u)) \frac{\partial (x_j, x_{i_1} \cdots x_{i_{k-1}})}{\partial (u_1, \dots, u_k)} du \end{aligned}$$

where the partial derivatives of $A_{1l}^I = (-1)^{1+l} \frac{\partial (x_{i_1}, \dots, x_{i_{k-1}})}{\partial (u_1, \dots, \hat{u}_l, \dots, u_k)}$ are defined in terms of Rademacher's theorem applied to the $k-1$ dimensional faces of $[a, b]$.

18.5 What Does it Mean?

Stokes theorem is a statement about integration by parts. However, one can give geometric meaning to what it says. These considerations will come from the area formula which I will use as needed.

For a particular l there are two faces in the boundary term for the Stokes formula. Consider the one where the l^{th} component is b_l . Recall that J was the set of increasing lists of $k-1$ indices.

$$\int_{[a,b]_l} \sum_{I \in J} ((\alpha_I \circ r) A_{1l}^I)(u_l(b_l)) du_l$$

Here $A_{1l}^I = (-1)^{1+l} \frac{\partial(x_{i_1}, \dots, x_{i_{k-1}})}{\partial(u_1, \dots, \hat{u}_l, \dots, u_k)}$ and $I = (i_1, \dots, i_{k-1})$. Letting

$$J_*(u_l) = \sqrt{\sum_{I \in J} (A_{1l}^I)^2(u_l(b_l))},$$

this term is of the form

$$\int_{[a,b]_l} \sum_{I \in J} \left((\alpha_I \circ r) \frac{A_{1l}^I}{J_*(u_l)} \right) (u_l(b_l)) J_*(u_l) du_l \quad (18.8)$$

Define $\frac{A_{1l}^I}{J_*(u_l)} = 0$ if $J_*(u_l) = 0$ on Z_l . By Lemma 17.3.1, $\mathcal{H}^{k-1}(r_l(Z_l)) = 0$ so the considerations presented here hold off a set of \mathcal{H}^{k-1} measure zero in $r_l([a, b]_l)$. Also we can ignore the set where the derivative does not exist thanks to Lemma 17.1.2 which says Lipschitz mapse of sets of measure zero have measure zero. Using the Binet Cauchy theorem to identify $J_*(u_l)$ with $(\det(Dr_{b_l}(u_l)^* Dr_{b_l}(u_l)))^{1/2}$, 18.8 reduces to

$$\int_{r_{b_l}([a,b]_l)} \#(x) \sum_{I \in J} \alpha_I(x) N_{b_l}^I(x) d\mathcal{H}^{k-1}(x)$$

where $N_{b_l}^I(r_{b_l}(u_l)) = \frac{A_{1l}^I}{J_*(u_l)}(u_l(b_l)) = \frac{1}{J_*(u_l)}(-1)^{1+l} \frac{\partial(x_{i_1}, \dots, x_{i_{k-1}})}{\partial(u_1, \dots, \hat{u}_l, \dots, u_k)}$ is a component of a unit vector in $\mathbb{R}^{C(p,k-1)}$ at least \mathcal{H}^{k-1} a.e. Assume that r_{b_l} is one to one or is one to one off a set S which has $\mathcal{H}^{k-1}(r_{b_l}(S)) = 0$. That way we can eliminate $\#(x)$ the number of times x is hit by r_{b_l} replacing it with 1. Thus, generalizing the notation, the boundary term in Stokes theorem is of the form

$$\begin{aligned} & \sum_{l=1}^k \int_{r_{b_l}([a,b]_l)} \sum_{I \in J} \alpha_I(x) N_{b_l}^I(x) d\mathcal{H}^{k-1}(x) \\ & - \sum_{l=1}^k \int_{r_{a_l}([a,b]_l)} \sum_{I \in J} \alpha_I(x) N_{a_l}^I(x) d\mathcal{H}^{k-1}(x) \end{aligned}$$

where $\sum_{I \in J} (N_{b_l}^I(x))^2 = 1$. Letting $N^I = N_{b_l}^I$ on $r_{b_l}([a, b]_l)$ and $-N_{a_l}^I$ on $r_{a_l}([a, b]_l)$, the above is of the form

$$\int_{r(\partial[a,b])} \sum_{I \in J} \alpha_I(x) N^I(x) d\mathcal{H}^{k-1}$$

Similarly

$$\begin{aligned} \int_r d\omega & \equiv \sum_{I \in J} \int_{[a,b]} \sum_{j=1}^p \frac{\partial \alpha_I}{\partial x_j}(r(u)) \frac{\partial(x_j, x_{i_1}, \dots, x_{i_{k-1}})}{\partial(u_1, \dots, u_k)} du \\ & = \int_{[a,b]} \sum_{j=1}^p \sum_{I \in J} \frac{\partial \alpha_I}{\partial x_j}(r(u)) \frac{\partial(x_j, x_{i_1}, \dots, x_{i_{k-1}})}{\partial(u_1, \dots, u_k)} du \end{aligned}$$

Now that determinant is only nonzero if j is none of the i_s . By the Binet Cauchy theorem,

$$\det(D\mathbf{r}(\mathbf{u})^* D\mathbf{r}(\mathbf{u})) = J_*(\mathbf{u})^2 = \sum_{x_{i_1} < \dots < x_{i_k}} \left(\frac{\partial(x_{i_1} \dots x_{i_k})}{\partial(u_1, \dots, u_k)} \right)^2$$

and so it follows from the area formula that there exists N_j^I for each I an increasing list of $k-1$ indices such that for J all such increasing lists, $\sum_{j=1}^p \sum_{I \in J} (N_j^I)^2 = 1$ and

$$\begin{aligned} \int_{\mathbf{r}} d\omega &= \int_{[a,b]} \sum_{j=1}^p \sum_{I \in J} \frac{\partial \alpha_I}{\partial x_j}(\mathbf{r}(\mathbf{u})) \frac{\frac{\partial(x_j, x_{i_1}, \dots, x_{i_{k-1}})}{\partial(u_1, \dots, u_k)}}{J_*(\mathbf{u})} J_*(\mathbf{u}) d\mathbf{u} \\ &= \int_{\mathbf{r}([a,b])} \#(\mathbf{x}) \sum_{I \in J} \sum_{j=1}^p \frac{\partial \alpha_I}{\partial x_j}(\mathbf{x}) N_j^I(\mathbf{x}) d\mathcal{H}^k \end{aligned} \quad (18.9)$$

where $\#(\mathbf{x})$ is the number of times \mathbf{x} is hit by \mathbf{r} . Thus if \mathbf{r} is one to one off \hat{S} where $\mathbf{r}(\hat{S})$ has \mathcal{H}^k measure zero, it follows that we can remove $\#(\mathbf{x})$ from the formula. As before, we can ignore the set where $J_*(\mathbf{u}) = 0$ thanks to Lemma 17.3.1. Also we can ignore the set of \mathbf{u} where the function is not differentiable by Lemma 17.1.2.

Observation 18.5.1 *Stokes theorem may be thought of as a statement about $\mathbf{r}([a,b])$ and $\mathbf{r}(\partial[a,b])$ which involves geometrical concepts dependent on these sets. This holds whenever Lipschitz \mathbf{r} restricted to each face of $[a,b]$ is one to one off a set S where $\mathbf{r}(S)$ has \mathcal{H}^{k-1} measure zero, and \mathbf{r} is one to one off \hat{S} where $\mathbf{r}(\hat{S})$ has \mathcal{H}^k measure zero.*

I think that the most important case is where $k = p$ and in this case we have the divergence theorem. Here there are exactly $p-1$ increasing lists of indices I and these are of the form $(1, \dots, \hat{j}, \dots, p)$. We let α_I be denoted as $a_j(-1)^{j+1}$ where j is the index missing in I . Therefore, the formula for $\int_{\mathbf{r}} d\omega$ reduces to

$$\pm \int_{\mathbf{r}([a,b])} \sum_{j=1}^p \frac{\partial a_j}{\partial x_j}(\mathbf{x}) d\mathcal{H}^k$$

assuming that \mathbf{r} is one to one off a set S for which $\mathbf{r}(S)$ has \mathcal{H}^p measure zero. This is because from 18.9

$$\frac{\partial \alpha_I}{\partial x_j}(\mathbf{r}(\mathbf{u})) \frac{\frac{\partial(x_j, x_{i_1}, \dots, x_{i_{k-1}})}{\partial(u_1, \dots, u_k)}}{J_*(\mathbf{u})} J_*(\mathbf{u}) = ((-1)^{j+1})^2 a_i(\mathbf{r}(\mathbf{u})) \frac{\frac{\partial(x_1, \dots, x_j, \dots, x_p)}{\partial(u_1, \dots, u_k)}}{J_*(\mathbf{u})} = \pm 1$$

The boundary terms reduce to

$$\sum_{l=1}^k \int_{\mathbf{r}_{b_l}([a,b]_l)} \sum_{j=1}^p a_j(\mathbf{x}) N^j(\mathbf{x}) d\mathcal{H}^{k-1}(\mathbf{x}) - \sum_{l=1}^k \int_{\mathbf{r}_{a_l}([a,b]_l)} \sum_{j=1}^p a_j(\mathbf{x}) N^j(\mathbf{x}) d\mathcal{H}^{k-1}(\mathbf{x})$$

where for $\mathbf{x} = \mathbf{r}(\mathbf{u})$, $N^j(\mathbf{x}) = (-1)^{j+1} \frac{\partial(x_1, \dots, \hat{x}_j, \dots, x_p)}{\partial(u_1, \dots, \hat{u}_j, \dots, u_p)}$. This can be written in the form

$$\pm \int_{\mathbf{r}([a,b])} \sum_{j=1}^p \frac{\partial a_j}{\partial x_j}(\mathbf{x}) d\mathcal{H}^k = \int_{\mathbf{r}(\partial[a,b])} \sum_{j=1}^p a_j(\mathbf{x}) N^j(\mathbf{x}) d\mathcal{H}^{k-1}(\mathbf{x})$$

This divergence theorem is discussed more later when also the unit vector N whose j^{th} component is N^j is described as an **exterior** unit normal provided $\frac{\partial(x_1, \dots, x_p)}{\partial(u_1, \dots, u_k)} \geq 0$. Note how the $(-1)^{j+1}$ on the a_j is responsible for the $(-1)^{j+l}$ instead of $(-1)^{1+l}$ in the description of N^j .

18.6 Examples of $r([a, b])$

I want to point out that there are many examples of $r([a, b])$ which fit into the above integration by parts idea of Stokes theorem in which r is Lipschitz, in order to tie this more to the way we usually think of these theorems in Calculus. Consider the following picture in which a closed ball B of radius 1 is inscribed into the box $Q \equiv [-1, 1] \equiv \prod_{i=1}^k [-1, 1]$. Let P be the projection map onto this ball B .



Then it is geometrically obvious that the projection map P satisfies $P(Q \setminus B) = \partial B$ a set of \mathcal{H}^k measure zero and that $P: \partial[-1, 1] \rightarrow \partial B$ is one to one and onto on $\partial[-1, 1]$. Now let $r: B \rightarrow \mathbb{R}^p$ for $p \geq k$ be Lipschitz and one to one. Then $r \circ P: [-1, 1] \rightarrow r(B)$ satisfies all the necessary conditions for an application of Stokes theorem including the geometric descriptions just given. What kinds of sets are in $r(B)$ for B a closed ball and r Lipschitz and one to one? I think you can see that this would include virtually everything of interest. You could stretch B in various directions, pinch it, bend it, etc. Roughly speaking, imagine a ball of soft clay and doing what a child would do to it before he tears it into little pieces, throws them around the room and stomps them into the carpet. The result would be one of the possible sets $r(B)$. Since r is a homeomorphism, the interior of B corresponds to the relative interior of $r(B)$, points $x \in r(B) = r \circ P([-1, 1])$ which are not in $r(\partial B)$. The boundary faces of $[-1, 1]$ and $r \circ P$ restricted to these faces will parametrize finitely many disjoint pieces of $r(\partial B)$.

Of course you could also consider chains of such boxes as described earlier in the case that r is C^1 . However, when you can allow r to be Lipschitz, it is clear that the theory is sufficiently general to include most things which would be of interest in any application from a single box. Next is a discussion of orientation placed here to make an analogy with the case of line integrals and oriented curves.

18.7 Orientation and Degree

Here I will consider orientation briefly. As in the case of a curve, it reduces to considerations of $r^{-1} \circ \hat{r}$.

Proposition 18.7.1 Suppose $r([a, b]) = \hat{r}([\hat{a}, \hat{b}])$, two sets in \mathbb{R}^p and both r, \hat{r} are one to one and Lipschitz, $[a, b], [\hat{a}, \hat{b}]$ being two parameter domains in $\mathbb{R}^k, k \leq p$. Then for

$$\omega = a(x) dx_{i_1} \wedge \cdots \wedge dx_{i_k}$$

Assume also that $r^{-1} \circ \hat{r}$ is Lipschitz and $\det(D(r^{-1} \circ \hat{r})(t)) \geq 0$ a.e. This is a statement about orientation. It follows then that $\int_r \omega = \int_{\hat{r}} \omega$.

Proof: Let $(a, b) \equiv \prod_{j=1}^k (a_j, b_j)$. Now from the area formula

$$\begin{aligned} \int_r \omega &\equiv \int_{[a, b]} a(r(u)) \frac{\partial(x_{i_1}, \dots, x_{i_k})}{\partial(u_1, \dots, u_k)}(u) du \\ &= \int_{[\hat{a}, \hat{b}]} a(r(r^{-1} \circ \hat{r})(t)) \frac{\partial(x_{i_1}, \dots, x_{i_k})}{\partial(u_1, \dots, u_k)}(r^{-1} \circ \hat{r}(t)) \det(D(r^{-1} \circ \hat{r})(t)) dt \\ &= \int_{[\hat{a}, \hat{b}]} a(r(r^{-1} \circ \hat{r})(t)) \frac{\partial(x_{i_1}, \dots, x_{i_k})}{\partial(t_1, \dots, t_k)}(t) dt \equiv \int_{\hat{r}} \omega \blacksquare \end{aligned}$$

An application of the area formula gives the following corollary. I will use $r^{-1} \circ \hat{r}$ to denote a Lipschitz function which is one to one off a set S which equals the Lipschitz function $r^{-1} \circ \hat{r}$ on S^C . In particular if S has measure 0 and $r^{-1} \circ \hat{r}$ is Lipschitz on S^C , then you could extend to a Lipschitz function which would map S to a set of measure zero, thus being in the situation of this corollary.

Corollary 18.7.2 Suppose $r^{-1} \circ \hat{r}$ is one to one off a set $S \subseteq [\hat{a}, \hat{b}]$ and that $r^{-1} \circ \hat{r}(S)$ has measure zero. Then the above would hold with no change.

Since this allows for Lipschitz functions, this is slightly more general than the usual situation from Calculus even in one dimension. However, more can be said. Orientation is really a statement about the degree of the map $r^{-1} \circ \hat{r}$, a concept which makes perfect sense without any direct reference to differentiability.

Recall that with a smooth curve C having points p, q and a one to one map to this curve, there are two ways to move over the curve, from p to q or from q to p . One defines equivalence classes on the continuous mappings r which map a closed interval to C . Two of these r, \hat{r} are equivalent if $r^{-1} \circ \hat{r}$ is increasing. It follows from the intermediate value theorem of Bolzano and a simple argument that this composition of maps is either increasing or decreasing. Thus, from the theorem about differentiation of monotone functions, $(r^{-1} \circ \hat{r})'$ is nonnegative a.e. exactly when the two parametrizations give the same orientation. In the above, this is determined by $\det(D(r^{-1} \circ \hat{r}))$. In addition, this reduces to a topological notion having to do with the degree. Instead of “increasing” we say that $d(r^{-1} \circ \hat{r}, \Omega, (a, b))$ is 1. The notion of “increasing” is not available.

Recall that from Proposition 15.6.7, Corollary 15.6.6, $d(r^{-1} \circ \hat{r}, (\hat{a}, \hat{b}), (a, b))$ is either 1 or -1 . This is the topological degree of the mapping $r^{-1} \circ \hat{r}$ which is constant on the connected component (a, b) of $\mathbb{R}^k \setminus r^{-1} \circ \hat{r}(\partial([\hat{a}, \hat{b}]))$. Suppose then that this degree $d(r^{-1} \circ \hat{r}, (\hat{a}, \hat{b}), y)$ is 1. Then from Proposition 17.7.4 on Page 469, it follows that whenever $f \in C_c(a, b)$

$$\int f(u) d(r^{-1} \circ \hat{r}, \Omega, u) du = \int_{(\hat{a}, \hat{b})} f(r^{-1} \circ \hat{r}(t)) \det(D(r^{-1} \circ \hat{r})(t)) dt$$

You could take an increasing sequence $f_n(u) \rightarrow \mathcal{X}_{(a, b)}(u)$ of the above sort. From the area formula,

$$\int f_n(u) du = \int_{(\hat{a}, \hat{b})} f(r^{-1} \circ \hat{r}(t)) |\det(D(r^{-1} \circ \hat{r})(t))| dt$$

and so

$$0 = \int_{(\hat{a}, \hat{b})} f_n(\mathbf{r}^{-1} \circ \hat{\mathbf{r}}(t)) (|\det(D(\mathbf{r}^{-1} \circ \hat{\mathbf{r}})(t))| - \det(D(\mathbf{r}^{-1} \circ \hat{\mathbf{r}})(t))) dt.$$

Using the monotone convergence theorem, it follows that

$$0 = \int_{(\hat{a}, \hat{b})} (|\det(D(\mathbf{r}^{-1} \circ \hat{\mathbf{r}})(t))| - \det(D(\mathbf{r}^{-1} \circ \hat{\mathbf{r}})(t))) dt.$$

and so $\det(D(\mathbf{r}^{-1} \circ \hat{\mathbf{r}})(t)) \geq 0$ a.e. This is the condition given in the above proposition. Thus this condition which applies to Lipschitz functions follows from a statement that $d(\mathbf{r}^{-1} \circ \hat{\mathbf{r}}, \Omega, \mathbf{u}) = 1$. Conversely, if the condition $\det(D(\mathbf{r}^{-1} \circ \hat{\mathbf{r}})(t)) \geq 0$ a.e. it will follow from the above formula that $d(\mathbf{r}^{-1} \circ \hat{\mathbf{r}}, \Omega, \mathbf{u}) = 1$ since it cannot equal -1 .

18.8 Examples of Stoke's Theorem

Here the attempt is made to tie this formalism to the usual things studied in calculus.

18.8.1 Fundamental Theorem of Calculus

First let $k = p = 1$. What does Stoke's theorem say? In this case, $d\omega$ is a 1 form and so ω should be a 0 form which is just a function. $\omega = a(x)$, $d\omega = a'(x)dx$. Then if $r : [a, b] \rightarrow \mathbb{R}$ is C^1 ,

$$\begin{aligned} \int_r d\omega &\equiv \int_r a'(x)dx \equiv \int_a^b a'(r(u)) \frac{dr}{du} du = a(r(b)) - a(r(a)) \\ \int_{\partial r} \omega &= \int_{\partial r} a(x) = (-1)^2 a(r(b)) - (-1)^2 a(r(a)) \end{aligned}$$

which is the same thing. It is just the fundamental theorem of calculus essentially.

18.8.2 Line Integrals for Conservative Fields

What if $p = 3, k = 1$? This is similar. You need to have ω a zero form. Thus $\omega = a(x_1, \dots, x_m)$. Then $d\omega = \sum_i a_{x_i} dx_i$. Letting $r : [a, b] \rightarrow \mathbb{R}^m$,

$$\begin{aligned} \int_r d\omega &= \int_a^b \nabla a \cdot \mathbf{r}' du = a(r(b)) - a(r(a)) \\ \int_{\partial r} \omega &= (-1)^{(1+1)} a(r(b)) - (-1)^{(1+1)} a(r(a)) \end{aligned}$$

which is the same thing. This is the case of a conservative vector field, the potential function being a .

18.8.3 Green's Theorem

Next ω be a 1 form and let $p = k = 2$ so $d\omega$ is a 2 form and ω is a 1 form. Say

$$\begin{aligned} \omega &= P(x, y) dx + Q(x, y) dy \\ d\omega &= P_y(x, y) dy \wedge dx + Q_x dx \wedge dy \end{aligned}$$

Recall that terms like $dx \wedge dx$ are zero because they result in a determinant which equals 0. Then let \mathbf{r} be smooth and map $[a, b] \times [c, d]$ to \mathbb{R}^2 .

$$\mathbf{r}(s, u) \equiv (x(s, u), y(s, u))^T.$$

Then $\int_{\mathbf{r}} d\omega \equiv$

$$\begin{aligned} & \int_a^b \int_c^d P_y(x(s, u), y(s, u)) \begin{vmatrix} y_s & y_u \\ x_s & x_u \end{vmatrix} + Q_x(x(s, u), y(s, u)) \begin{vmatrix} x_s & x_u \\ y_s & y_u \end{vmatrix} ds du \\ &= \int_a^b \int_c^d [Q_x(x(s, u), y(s, u)) - P_y(x(s, u), y(s, u))] (x_s y_u - x_u y_s) ds du \end{aligned}$$

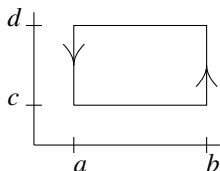
Let $U = \mathbf{r}([a, b] \times [c, d])$. Then by the change of variables theorem, this equals

$$\int_U (Q_x(x, y) - P_y(x, y)) \operatorname{sgn}(x_s y_u - x_u y_s) dx dy$$

Next consider $\int_{\partial \mathbf{r}} \omega$. For $\mathbf{r}(s, u) = (x(s, u), y(s, u))^T$, this equals

$$\begin{aligned} & \int_c^d P(x(1, u), y(1, u)) \frac{\partial x(1, u)}{\partial u} + Q(x(1, u), y(1, u)) \frac{\partial y(1, u)}{\partial u} du \\ & - \int_c^d P(x(0, u), y(0, u)) \frac{\partial x(0, u)}{\partial u} + Q(x(0, u), y(0, u)) \frac{\partial y(0, u)}{\partial u} du \\ & + \int_a^b P(x(s, 0), y(s, 0)) \frac{\partial x(s, 0)}{\partial s} + Q(x(s, 0), y(s, 0)) \frac{\partial y(s, 0)}{\partial s} ds \\ & - \int_a^b P(x(s, 1), y(s, 1)) \frac{\partial x(s, 1)}{\partial s} + Q(x(s, 1), y(s, 1)) \frac{\partial y(s, 1)}{\partial s} ds \end{aligned}$$

This is computing a line integral by summing the contributions around the edges of $[a, b] \times [c, d]$. It is $\int_C P dx + Q dy$ where the orientation on this curve C comes from the counter clockwise orientation on the boundary of $[a, b] \times [c, d]$ as shown in the picture.



Thus if $x_s y_u - x_u y_s > 0$, this is just Green's theorem from calculus. Thus this gives a proof of this important theorem provided the region U is the C^1 image of a rectangle. One could generalize to consider chains of rectangles which, as mentioned above yields fairly general surfaces in \mathbb{R}^2 . One could also include Lipschitz maps for even more generality. However, this falls short of the best results for Green's theorem which involve a rectifiable simple closed curve with U its interior. Rectifiable only requires the curve to have finite length.

18.8.4 Stoke's Theorem from Calculus

Next let $k = 2$ and $p = 3$. Thus $d\omega$ is a 2 form so ω is a 1 form. Say

$$\omega = P(x, y, z) dx + Q(x, y, z) dy + R(x, y, z) dz$$

Then

$$d\omega = P_y dy \wedge dx + P_z dz \wedge dx + Q_x dx \wedge dy + Q_z dz \wedge dy + R_x dx \wedge dz + R_y dy \wedge dz$$

Now let $\mathbf{r} : [0, 1]^2 \rightarrow \mathbb{R}^3$, $\mathbf{r}(s, u) = (x(s, u), y(s, u), z(s, u))$. Then letting the various functions be defined at $\mathbf{r}(s, u)$,

$$\begin{aligned} \int_{\mathbf{r}} d\omega = \int_{[0,1]^2} & \left(P_y \det \begin{pmatrix} y_s & y_u \\ x_s & x_u \end{pmatrix} + P_z \det \begin{pmatrix} z_s & z_u \\ x_s & x_u \end{pmatrix} + Q_x \det \begin{pmatrix} x_s & x_u \\ y_s & y_u \end{pmatrix} \right. \\ & \left. + Q_z \det \begin{pmatrix} z_s & z_u \\ y_s & y_u \end{pmatrix} + R_x \det \begin{pmatrix} x_s & x_u \\ z_s & z_u \end{pmatrix} + R_y \det \begin{pmatrix} y_s & y_u \\ z_s & z_u \end{pmatrix} \right) ds du \end{aligned}$$

This equals

$$\begin{aligned} & \int_{[0,1]^2} (R_y - Q_z)(y_s z_u - y_u z_s) + (P_z - R_x)(x_u z_s - x_s z_u) \\ & + (Q_x - P_y)(x_s y_u - x_u y_s) ds du \end{aligned}$$

By the definition of surface area on $S \equiv \mathbf{r}([0, 1]^2)$, see Definition 14.2.1 the area increment on the surface $\mathbf{r}([0, 1]^2)$ is

$$\begin{aligned} & \sqrt{|y_s z_u - y_u z_s|^2 + |x_u z_s - x_s z_u|^2 + |x_s y_u - x_u y_s|^2} ds du \\ & = \det(D\mathbf{r}(u, s))^* D\mathbf{r}(s, u))^{1/2} ds du \end{aligned}$$

also the vector $(R_y - Q_z)\mathbf{i} + (P_z - R_x)\mathbf{j} + (Q_x - P_y)\mathbf{k}$ is the curl of the vector field $\mathbf{F} \equiv P\mathbf{i} + Q\mathbf{j} + R\mathbf{k}$. Thus in more familiar calculus notation, the above integral is of the form $\int_{S \equiv \mathbf{r}([0,1]^2)} \nabla \times \mathbf{F} \cdot \mathbf{p} dS$ where \mathbf{p} is a vector which is in the direction of

$$(y_s z_u - y_u z_s)\mathbf{i} + (x_u z_s - x_s z_u)\mathbf{j} + (x_s y_u - x_u y_s)\mathbf{k}.$$

It happens that this vector \mathbf{p} is perpendicular to the surface S at every point where $\mathbf{r}_s \times \mathbf{r}_u \neq \mathbf{0}$, which is seen to occur in the above formula. Recall that this follows from noting that, as (hopefully) discussed in beginning calculus, you have $\mathbf{r}_s \cdot \mathbf{r}_s \times \mathbf{r}_u = \mathbf{r}_u \cdot \mathbf{r}_s \times \mathbf{r}_u = 0$ and this is sufficient to claim, based on geometric reasoning that it is perpendicular to the surface. Thus $\int_{\mathbf{r}} d\omega$ is one side of the usual Stoke's theorem from calculus. You could generalize to chains of rectangles as well.

Next consider what happens with $\int_{\partial \mathbf{r}} \omega$. This is just like it was with Green's theorem but with more terms. To save on space, $P(x(1, u), y(1, u), z(1, u))$ is denoted as $P(1, u)$ with similar considerations for Q and R . Then this results in

$$\begin{aligned} & \int_0^1 \left(P(1, u) \frac{\partial x(1, u)}{\partial u} + Q(1, u) \frac{\partial y(1, u)}{\partial u} + R(1, u) \frac{\partial z(1, u)}{\partial u} \right) du \\ & - \int_0^1 \left(P(0, u) \frac{\partial x(0, u)}{\partial u} + Q(0, u) \frac{\partial y(0, u)}{\partial u} + R(1, u) \frac{\partial z(1, u)}{\partial u} \right) du \\ & + \int_0^1 P(s, 0) \frac{\partial x(s, 0)}{\partial s} + Q(s, 0) \frac{\partial y(s, 0)}{\partial s} + R(s, 0) \frac{\partial z(s, 0)}{\partial s} ds \\ & - \int_0^1 P(s, 1) \frac{\partial x(s, 1)}{\partial s} + Q(s, 1) \frac{\partial y(s, 1)}{\partial s} + R(s, 1) \frac{\partial z(s, 1)}{\partial s} ds \end{aligned}$$

As before, this is a line integral of the form $\int_C Pdx + Qdy + Rdz$, where C is an oriented curve bounding the surface S . This orientation will determine the direction of the vector \mathbf{p} above.

18.8.5 The Divergence Theorem

In this case, we have a parameter domain $[\mathbf{a}, \mathbf{b}] \subseteq \mathbb{R}^p$ and the differential form is

$$\omega = \sum_{r=1}^p \alpha_r(\mathbf{x}) (-1)^{r-1} dx_1 \wedge \cdots \wedge d\hat{x}_r \wedge \cdots \wedge dx_p$$

where $d\hat{x}_r$ with the hat means that dx_r is omitted. The reason for the $(-1)^{r-1}$ is to make minus signs disappear in $d\omega$. This led to the divergence theorem

$$\pm \int_{\mathbf{r}([\mathbf{a}, \mathbf{b}])} \sum_{j=1}^p \frac{\partial \alpha_j}{\partial x_j}(\mathbf{x}) d\mathcal{H}^k = \int_{\mathbf{r}(\partial[\mathbf{a}, \mathbf{b}])} \sum_{j=1}^p \alpha_j(\mathbf{x}) N^j(\mathbf{x}) d\mathcal{H}^{k-1}(\mathbf{x})$$

where the $+$ sign holds if and only if $\frac{\partial(x_1, \dots, x_p)}{\partial(u_1, \dots, u_p)} \geq 0$. We needed to have \mathbf{r} is one to one off a set \hat{S} where $\mathbf{r}(\hat{S})$ has \mathcal{H}^p measure zero and $\mathbf{r}_{b_l}, \mathbf{r}_{a_l}$ are also one to one or at least one to one off a set S where $\mathbf{r}_{b_l}(S), \mathbf{r}_{a_l}(S)$ have \mathcal{H}^{p-1} measure zero. It remains to consider the vector \mathbf{N} which has j^{th} component N^j . I want to argue that this vector is a.e. normal to $\mathbf{r}(\partial[\mathbf{a}, \mathbf{b}])$ and that sometimes it is an outer normal. Recall that for $\mathbf{x} = \mathbf{r}_{b_l}(\mathbf{u}_l)$, $N^j(\mathbf{x}) = (-1)^{j+l} \frac{\partial(x_1, \dots, \hat{x}_j, \dots, x_p)}{\partial(u_1, \dots, \hat{u}_l, \dots, u_p)} \frac{1}{J_*(\mathbf{x})}$. Thus for a point on the boundary, for $i \neq l$,

$$\begin{aligned} \mathbf{x}_{u_i} \cdot \mathbf{N} &= \sum_{j=1}^p x_{j, u_i} N^j(\mathbf{x}) = \sum_{j=1}^p x_{j, u_i} (-1)^{l+j} \frac{\partial(x_1, \dots, \hat{x}_j, \dots, x_p)}{\partial(u_1, \dots, \hat{u}_l, \dots, u_p)} \frac{1}{J_*(\mathbf{x})} \\ &= \frac{1}{J_*(\mathbf{x})} \det \begin{pmatrix} \mathbf{x}_{u_1} & \cdots & \mathbf{x}_{u_i} & \cdots & \mathbf{x}_{u_i} & \cdots & \mathbf{x}_{u_p} \end{pmatrix} = 0 \end{aligned}$$

since it is a determinant with two equal columns. This involved expanding along the l^{th} column which was filled by \mathbf{x}_{u_i} . Also, expanding along this column,

$$\mathbf{x}_{u_l} \cdot \mathbf{N} = \frac{1}{J_*(\mathbf{x})} \det(D\mathbf{r}) = \text{sgn}(\det(D\mathbf{r}))$$

It follows that \mathbf{N} is perpendicular to $\mathbf{r}(\partial[\mathbf{a}, \mathbf{b}])$ and that the angle between \mathbf{x}_{u_l} and \mathbf{N} is no more than 90 degrees if $\frac{\partial(x_1, \dots, x_p)}{\partial(u_1, \dots, u_p)} \geq 0$. Now on the face where $u_l = b_l$, \mathbf{x}_{u_l} points away from this face and so \mathbf{N} points in roughly the same direction and is an “**outer**” normal. Similar considerations apply when $u_l = a_l$ but here the $-\mathbf{x}_{u_l}$ points away and we use $-\mathbf{N}$ because of the subtraction in the boundary integrals. This yields the following.

Theorem 18.8.1 *Let \mathbf{r} be Lipschitz and one to one off S where $\mathbf{r}(S)$ has \mathcal{H}^p measure zero and suppose a similar condition holds for \mathbf{r}_{a_l} and \mathbf{r}_{b_l} . Let α_j be C^1*

$$\omega = \sum_{j=1}^p \alpha_j(\mathbf{x}) (-1)^{j-1} dx_1 \wedge \cdots \wedge d\hat{x}_j \wedge \cdots \wedge dx_p, \quad \alpha(\mathbf{x}) = (\alpha_1(\mathbf{x}), \dots, \alpha_p(\mathbf{x}))$$

and suppose $\det(D\mathbf{r}(\mathbf{u})) \geq 0$. Then

$$\int_{\mathbf{r}} d\omega = \int_{\mathbf{r}([a,b])} \sum_j \frac{\partial \alpha_j}{\partial x_j}(\mathbf{x}) d\mathcal{H}^p = \int_{\mathbf{r}(\partial[a,b])} \boldsymbol{\alpha} \cdot \mathbf{N} d\mathcal{H}^{p-1} \quad (18.10)$$

where \mathbf{N} is an outer unit normal in the sense that the angle between the vector \mathbf{x}_{u_i} and \mathbf{N} is no more than 90 degrees if $\frac{\partial(x_1, \dots, x_p)}{\partial(u_1, \dots, u_p)} \geq 0$.

Also, if you know the divergence theorem, then you can directly give the usual Calculus version of Green's and Stoke's theorems from Calculus. This is developed in the exercises.

18.9 The Reynolds Transport Formula

The Reynolds transport formula is a generalization of the formula for taking the derivative under an integral. It depends on the divergence theorem. I will use the chain rule of Theorem 17.3.5 as needed.

$$\frac{d}{dt} \int_{a(t)}^{b(t)} f(x, t) dx = \int_{a(t)}^{b(t)} \frac{\partial f}{\partial t}(x, t) dx + f(b(t), t) b'(t) - f(a(t), t) a'(t)$$

First is an interesting lemma about the determinant. A $p \times p$ matrix can be thought of as a vector in \mathbb{C}^{p^2} . Just imagine stringing it out into one long list of numbers. In fact, a way to give the norm of a matrix is just $\sum_i \sum_j |A_{ij}|^2 \equiv \|A\|^2$. This is called the Frobenius norm for a matrix. It makes no difference since all norms are equivalent, but this one is convenient in what follows. Also recall that \det maps $p \times p$ matrices to \mathbb{C} . It makes sense to ask for the derivative of \det on the set of invertible matrices, an open subset of \mathbb{C}^{p^2} with the norm measured as just described because $A \rightarrow \det(A)$ is continuous, so the set where $\det(A) \neq 0$ would be an open set. Recall from linear algebra that the sum of the entries on the main diagonal satisfies $\text{trace}(AB) = \text{trace}(BA)$ whenever both products make sense. Indeed, $\text{trace}(AB) \equiv \sum_i \sum_j A_{ij} B_{ji} = \text{trace}(BA)$

This next lemma is a very interesting observation about the determinant of a matrix added to the identity.

Lemma 18.9.1 $\det(I + U) = 1 + \text{trace}(U) + o(U)$ where $o(U)$ is defined in terms of the Frobenius norm for $p \times p$ matrices.

Proof:

$$\det(I + U) \equiv \sum_{i_1, i_2, \dots, i_p} \text{sgn}(i_1, i_2, \dots, i_p) (\delta_{i_1 1} + U_{i_1 1}) \cdots (\delta_{i_p p} + U_{i_p p})$$

which equals $\det(I)$ added to $\text{trace}(U)$ added to a sum of higher order terms of products of the U_{ij} . The $\text{trace}(U)$ comes from using only one U_{ij} in the above product. The resulting term will be 0 unless $i = j$ and so the end result of these will be $\text{trace}(U)$. Of course if more of the U_{ij} are included in the product, this yields the $o(U)$ term. ■

Of course, by equivalence of norms, one could use any other norm for the $p \times p$ matrices.

With this lemma, it is easy to find $D \det(F)$ whenever F is invertible.

$$\begin{aligned} \det(F + U) &= \det(F(I + F^{-1}U)) = \det(F) \det(I + F^{-1}U) \\ &= \det(F) (1 + \text{trace}(F^{-1}U) + o(U)) \\ &= \det(F) + \det(F) \text{trace}(F^{-1}U) + o(U) \end{aligned}$$

Therefore, $\det(F + U) - \det(F) = \det(F) \text{trace}(F^{-1}U) + o(U)$. This proves the following.

Proposition 18.9.2 *Let F^{-1} exist. Then $D \det(F)(U) = \det(F) \text{trace}(F^{-1}U)$.*

From this, suppose $F(t)$ is a $p \times p$ matrix and all entries are differentiable. Then the following describes $\frac{d}{dt} \det(F)(t)$.

Proposition 18.9.3 *Let $F(t)$ be a $p \times p$ matrix and all entries are at least Lipschitz. Then for a.e. t*

$$\frac{d}{dt} \det(F)(t) = \det(F(t)) \text{trace}(F^{-1}(t) F'(t)) = \det(F(t)) \text{trace}(F'(t) F^{-1}(t)) \quad (18.11)$$

Proof: From the above,

$$\begin{aligned} &\det(F(t+h)) - \det(F(t)) \\ &= \det(F(t)) \text{trace}(F^{-1}(F(t+h) - F(t))) + o(F(t+h) - F(t)) \end{aligned}$$

$\stackrel{=o(h) \text{ since } F' \text{ exists}}{\quad}$

Dividing by h and taking a limit yields 18.11. ■

Let $\mathbf{y} = \mathbf{h}(t, \mathbf{x})$ with $F = F(t, \mathbf{x}) = D_2 \mathbf{h}(t, \mathbf{x})$. I will write $\nabla_{\mathbf{y}}$ to indicate the gradient with respect to the \mathbf{y} variables and F' to indicate $\frac{\partial}{\partial t} F(t, \mathbf{x})$. I will be assuming what is needed to use the various theorems. In particular let \mathbf{h} be differentiable and one to one in \mathbf{x} . Note that $\mathbf{h}(t, \mathbf{x}) = \mathbf{y}$ and so by the inverse function theorem, or actually Corollary 8.10.6, this defines \mathbf{x} as a function of \mathbf{y} , also differentiable as \mathbf{h} because it is always assumed $\det F > 0$.

Now let V_t be $\mathbf{h}(t, V_0)$ where V_0 is an open bounded set. Let V_0 have a Lipschitz boundary so one can use the divergence theorem on V_0 . Thus this is concerned with smooth motion of a bounded open set with Lipschitz boundary. Let $(t, \mathbf{y}) \rightarrow \mathbf{f}(t, \mathbf{y})$ be Lipschitz. The idea is to simplify $\frac{d}{dt} \int_{V_t} \mathbf{f}(t, \mathbf{y}) dm_p(\mathbf{y})$. This will involve the change of variables in which the Jacobian will be $\det(F)$ which is assumed positive thus preserving the orientation of the normal vector for V_0 and V_t . In applications of this theory, $\det(F) \leq 0$ is not physically possible. Since $\mathbf{h}(t, \cdot)$ is better than Lipschitz and the boundary of V_0 is Lipschitz, V_t will be such that one can use the divergence theorem because the composition of Lipschitz functions is Lipschitz. See Corollary 14.3.6. Then, using the dominated convergence theorem as needed along with the area formula,

$$\begin{aligned} \frac{d}{dt} \int_{V_t} \mathbf{f}(t, \mathbf{y}) dm_p(\mathbf{y}) &= \frac{d}{dt} \int_{V_0} \mathbf{f}(t, \mathbf{h}(t, \mathbf{x})) \det(F) dm_p(\mathbf{x}) \quad (18.12) \\ &= \int_{V_0} \frac{\partial}{\partial t} \mathbf{f}(\cdot, \mathbf{h}(\cdot, \mathbf{x})) \det(F) dm_p(\mathbf{x}) + \int_{V_0} \mathbf{f}(t, \mathbf{h}(t, \mathbf{x})) \frac{\partial}{\partial t} (\det(F)) dm_p(\mathbf{x}) \end{aligned}$$

$$\begin{aligned}
&= \int_{V_0} \frac{\partial}{\partial t} (\mathbf{f}(t, \mathbf{h}(t, \mathbf{x}))) \det(F) dm_p(\mathbf{x}) + \int_{V_0} \mathbf{f}(t, \mathbf{h}(t, \mathbf{x})) \operatorname{trace}(F'F^{-1}) \det(F) dm_p(\mathbf{x}) \\
&= \int_{V_0} \left(\frac{\partial}{\partial t} \mathbf{f}(t, \mathbf{h}(t, \mathbf{x})) + \sum_i \frac{\partial \mathbf{f}}{\partial y_i} \frac{\partial y_i}{\partial t} \right) \det(F) dm_p(\mathbf{x}) \\
&\quad + \int_{V_0} \mathbf{f}(t, \mathbf{h}(t, \mathbf{x})) \operatorname{trace}(F'F^{-1}) \det(F) dm_p(\mathbf{x}) \\
&= \int_{V_t} \frac{\partial}{\partial t} \mathbf{f}(t, \mathbf{y}) dm_p(\mathbf{y}) + \int_{V_t} \sum_i \frac{\partial \mathbf{f}}{\partial y_i} \frac{\partial y_i}{\partial t} + \mathbf{f}(t, \mathbf{y}) \operatorname{trace}(F'F^{-1}) dm_p(\mathbf{y})
\end{aligned}$$

Now $\mathbf{v} \equiv \frac{\partial}{\partial t} \mathbf{h}(t, \mathbf{x})$ and also, as noted above, $\mathbf{y} \equiv \mathbf{h}(t, \mathbf{x})$ defines \mathbf{y} as a function of \mathbf{x} and so $\operatorname{trace}(F'F^{-1}) = \sum_{\alpha} \frac{\partial v_i}{\partial x_{\alpha}} \frac{\partial x_{\alpha}}{\partial y_i}$. Hence the double sum $\sum_{\alpha, i} \frac{\partial v_i}{\partial x_{\alpha}} \frac{\partial x_{\alpha}}{\partial y_i}$ is $\frac{\partial v_i}{\partial y_i} = \nabla_{\mathbf{y}} \cdot \mathbf{v}$. The above then gives

$$\begin{aligned}
&\int_{V_t} \frac{\partial}{\partial t} \mathbf{f}(t, \mathbf{y}) dm_p(\mathbf{y}) + \int_{V_t} \left(\sum_i \frac{\partial \mathbf{f}}{\partial y_i} \frac{\partial y_i}{\partial t} + \mathbf{f}(t, \mathbf{y}) \nabla_{\mathbf{y}} \cdot \mathbf{v} \right) dm_p(\mathbf{y}) \\
&= \int_{V_t} \frac{\partial}{\partial t} \mathbf{f}(t, \mathbf{y}) dm_p(\mathbf{y}) + \int_{V_t} (D_2 \mathbf{f}(t, \mathbf{y}) \mathbf{v} + \mathbf{f}(t, \mathbf{y}) \nabla_{\mathbf{y}} \cdot \mathbf{v}) dm_p(\mathbf{y}) \quad (18.13)
\end{aligned}$$

Now consider the i^{th} component of the second integral in the above. It is

$$\int_{V_t} \nabla_{\mathbf{y}} f_i(t, \mathbf{y}) \cdot \mathbf{v} + f_i(t, \mathbf{y}) \nabla_{\mathbf{y}} \cdot \mathbf{v} dm_p(\mathbf{y}) = \int_{V_t} \nabla_{\mathbf{y}} \cdot (f_i(t, \mathbf{y}) \mathbf{v}) dm_p(\mathbf{y})$$

At this point, use the divergence theorem to get this equals $\int_{\partial V_t} f_i(t, \mathbf{y}) \mathbf{v} \cdot \mathbf{n} d\mathcal{H}^{p-1}$. Therefore, from 18.13 and 18.12,

$$\frac{d}{dt} \int_{V_t} \mathbf{f}(t, \mathbf{y}) dm_p(\mathbf{y}) = \int_{V_t} \frac{\partial}{\partial t} \mathbf{f}(t, \mathbf{y}) dm_p(\mathbf{y}) + \int_{\partial V_t} \mathbf{f}(t, \mathbf{y}) \mathbf{v} \cdot \mathbf{n} d\mathcal{H}^{p-1} \quad (18.14)$$

this is the Reynolds transport formula.

Proposition 18.9.4 *Let $\mathbf{y} = \mathbf{h}(t, \mathbf{x})$ where \mathbf{h} is C^1 and let \mathbf{f} be Lipschitz continuous and let $V_t \equiv \mathbf{h}(t, V_0)$ where V_0 is a bounded open set which is on one side of a Lipschitz boundary so that the divergence theorem holds for V_0 . Then 18.14 is obtained.*

As with the divergence theorem, Some generalization should be possible to the case where \mathbf{h} giving the motion is only Lipschitz by using the version of the chain rule in Theorem 17.3.5 in the above argument when needed.

18.10 Exercises

1. Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ be given by $f(\mathbf{x}) \equiv \left(\sum_{i=1}^n \frac{x_i^2}{a_i^2} \right)^{1/2}$ where each $a_i > 0$. Thus for $y > 0$ $f^{-1}(y)$ is the boundary of an n dimensional ellipsoid. Using change of variables formula and the coarea formula, find the area of $f^{-1}(r)$.

2. Let $\pi : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be defined as $\pi(\mathbf{x}) = \mathbf{x}_i$ where $i = (i_1, i_2, \dots, i_m)$. What does Theorem 17.6.1 say if $\pi = \mathbf{f}$ in this theorem?
3. In calculus, you found the area of the parallelogram determined by two vectors \mathbf{u}, \mathbf{v} in \mathbb{R}^3 by taking the magnitude $|\mathbf{u} \times \mathbf{v}|$, meaning the Euclidean norm of the cross product. Show that you get the same answer by forming

$$(\det(\begin{pmatrix} \mathbf{u} & \mathbf{v} \end{pmatrix}))^{1/2}$$

where here you have $\begin{pmatrix} \mathbf{u} & \mathbf{v} \end{pmatrix}$ is the matrix which has \mathbf{u} as first column and \mathbf{v} as second column.

4. In calculus, you also found the volume of a parallelepiped determined by the vectors $\mathbf{u}, \mathbf{v}, \mathbf{w}$ by $|\mathbf{u} \cdot (\mathbf{v} \times \mathbf{w})|$, the absolute value of the box product. Show this turns out to be the same thing as

$$(\det(\begin{pmatrix} \mathbf{u} & \mathbf{v} & \mathbf{w} \end{pmatrix}))^{1/2}.$$

5. Imagine a fluid which does not move. Let $B(\mathbf{x}, \varepsilon)$ be a small ball in this fluid. Use the Euclidean norm. Then the force exerted on the ball of fluid is $-\int_{\partial B(\mathbf{x}, \varepsilon)} p \mathbf{n} dA$ where p is the pressure. Here \mathbf{n} is the unit exterior normal. Now the force acting on the ball from gravity is $-g \mathbf{k} \int_{B(\mathbf{x}, \varepsilon)} \rho dV$ where ρ signifies the density of the fluid and \mathbf{k} signifies the direction which is up. The vectors \mathbf{i}, \mathbf{j} are in the direction of the positive x and y axes respectively. These two forces add to $\mathbf{0}$ because it is given that the fluid does not move. Use the divergence theorem to show that $\nabla p = \rho g \mathbf{k}$. This is a really neat result.
6. Archimedes principle states that when a solid body is immersed in a static fluid, the force acting on the body by the fluid is directly up and equals the total weight of the fluid displaced. Surely this is an amazing result. It doesn't matter about the shape of the body. Remember that the weight is the acceleration of gravity times the mass. **Hint:** You need to start with the force acting on the body B by the fluid which is $-\int_{\partial B} p \mathbf{n} dA$. Assume the divergence theorem holds for B . As shown, this is typically the case.
7. You have a closed region R which is **fixed in space**. A fluid is flowing through R . The density of this fluid is $\rho(t, \mathbf{x})$ where \mathbf{x} gives the coordinates in space and ρ depends on t because it might change as time progresses. The velocity of this fluid is $\mathbf{v}(t, \mathbf{x})$. Then the rate at which the fluid crosses a surface S from one side to the other is $\int_S \rho \mathbf{v} \cdot \mathbf{n} dS$ where \mathbf{n} is the unit normal to the surface which points in the direction of interest. You can think about this a little and see that it is a reasonable claim. If, for example, $\mathbf{v} \cdot \mathbf{n} = 0$, then the velocity of the fluid would be parallel to the surface so it would not cross it at all. Also, the total mass of the fluid which is in R is $\int_R \rho dV$. Assuming anything you like about regularity, which is what we do in situations like this, explain using the divergence theorem why $\int_V \frac{\partial \rho}{\partial t} + \nabla \cdot (\rho \mathbf{v}) dV = 0$. Recall that $\nabla \cdot \mathbf{F}$ denotes the divergence of \mathbf{F} . Now explain why it should be the case that $\frac{\partial \rho}{\partial t} + \nabla \cdot (\rho \mathbf{v}) = 0$. This is called the balance of mass equation.
8. The permutation symbol is $\varepsilon_{ijk} = 1$ if the permutation $\begin{pmatrix} 1 & 2 & 3 \\ i & j & k \end{pmatrix}$ is even and -1 if this permutation is odd. Such a permutation is odd or even depending on

whether it requires an odd or even number of switches to obtain $(1, 2, 3)$. Thus $\varepsilon_{123} = 1$, $\varepsilon_{213} = -1$ and so forth. Show that $\sum_k \varepsilon_{ijk} \varepsilon_{irs} = \delta_{jr} \delta_{ks} - \delta_{js} \delta_{kr}$. We often agree to add over a repeated index to avoid having to write the summation symbol. Thus $\sum_k \varepsilon_{ijk} \varepsilon_{irs} = \varepsilon_{ijk} \varepsilon_{irs}$. If you do this, avoid having the repeated index repeated more than once in any term. Here δ_{ij} is 1 if $i = j$ and 0 if $i \neq j$.

9. Let U be a bounded open set in \mathbb{R}^p and suppose $u \in C^2(U) \cap C(\bar{U})$ such that $\nabla^2 u \geq 0$ in U . Then letting $\partial U = \bar{U} \setminus U$, it follows that

$$\max \{u(\mathbf{x}) : \mathbf{x} \in \bar{U}\} = \max \{u(\mathbf{x}) : \mathbf{x} \in \partial U\}.$$

The symbol ∇^2 is the Laplacian. Thus $\nabla^2 u = \sum_i u_{x_i x_i}$. In terms of repeated index summation convention, $\nabla^2 u = u_{,ii}$. **Hint:** Suppose this does not happen. Then there exists $\mathbf{x}_0 \in U$ with

$$u(\mathbf{x}_0) > \max \{u(\mathbf{x}) : \mathbf{x} \in \partial U\}.$$

Since U is bounded, there exists $\varepsilon > 0$ such that

$$u(\mathbf{x}_0) > \max \{u(\mathbf{x}) + \varepsilon |\mathbf{x}|^2 : \mathbf{x} \in \partial U\}.$$

Therefore, $u(\mathbf{x}) + \varepsilon |\mathbf{x}|^2$ also has its maximum in U because for ε small enough,

$$u(\mathbf{x}_0) + \varepsilon |\mathbf{x}_0|^2 > u(\mathbf{x}_0) > \max \{u(\mathbf{x}) + \varepsilon |\mathbf{x}|^2 : \mathbf{x} \in \partial U\}$$

for all $\mathbf{x} \in \partial U$. Now let \mathbf{x}_1 be the point in U where $u(\mathbf{x}) + \varepsilon |\mathbf{x}|^2$ achieves its maximum. Now recall the second derivative test from single variable calculus. Explain why at a local maximum of f you must have $\nabla^2 f \leq 0$. Apply this to the function $\mathbf{x} \rightarrow u(\mathbf{x}) + \varepsilon |\mathbf{x}|^2$ at the point \mathbf{x}_1 and get a contradiction. This is called the weak maximum principle.

10. Review the cross product from calculus. Show that in \mathbb{R}^3 , $(\mathbf{a} \times \mathbf{b})_i = \varepsilon_{ijk} a_j b_k$ where summation is over repeated indices. Using the above reduction identity of Problem 8, simplify $(\mathbf{a} \times \mathbf{b}) \times \mathbf{c}$ in terms of dot products. Then do the same for $\mathbf{a} \times (\mathbf{b} \times \mathbf{c})$.
11. Show that $\nabla \cdot (\nabla \times \mathbf{v}) = 0$. Now show $\int_{\partial V} \nabla \times \mathbf{v} \cdot \mathbf{n} dA = 0$ where V is a region for which the divergence theorem holds and \mathbf{v} is a C^2 vector field. $\nabla \times \mathbf{v}$ is the curl of \mathbf{v} . In the new notation, $(\nabla \times \mathbf{v})_i = \varepsilon_{ijk} \partial_j v_k$ where ∂_j is an operator which means to take the partial derivative with respect to x_j . It is understood here that the coordinates are rectangular coordinates. The first part of this is real easy if you remember the big theorem about equality of mixed partial derivatives.
12. Let U be a bounded open set in \mathbb{R}^2 which has a Lipschitz boundary so that the divergence theorem holds. Let the axes be oriented in the usual way as in calculus. Let $P(x, y), Q(x, y)$ be two smooth functions defined on \bar{U} . What does the divergence theorem say for the vector field $(Q, -P)$? If $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^2$ is Lipschitz with $\mathbf{r}(a) = \mathbf{r}(b)$ and \mathbf{r} is one to one with $\partial U = \mathbf{r}([a, b])$ and $\mathbf{r}'(t)$ is in the direction of $\mathbf{k} \times \mathbf{n}$ so $\mathbf{r}'(t) \times \mathbf{k}$ is in the direction of \mathbf{n} , a statement about orientation. Thus for $\mathbf{r}(t) = (x(t), y(t))$,

$$\mathbf{n} = c \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ x'(t) & y'(t) & 0 \\ 0 & 0 & 1 \end{vmatrix} = c(y'(t)\mathbf{i} - x'(t)\mathbf{j}), \text{ so } c = \frac{1}{\sqrt{y'(t)^2 + x'(t)^2}}$$

for some c a positive constant. (This is the cross product from calculus. Review this.), use the area formula to describe the boundary integral from the divergence theorem as an integral over $[a, b]$. We write this integral as $\int_{\partial U} Pdx + Qdy$. This yields a version of Green's theorem, $\int_U Q_x - P_y dm_2 = \int_{\partial U} Pdx + Qdy$ provided ∂U is a Lipschitz curve oriented as just described. However, more generality is possible, although I am not sure how far this has been generalized. You really only need to have the boundary of U be a rectifiable simple closed curve meaning it has finite length. In this setting, the Jordan curve theorem makes it possible to make sense of Green's theorem and in fact it holds. In what was just discussed, \bar{U} was the Lipschitz image of some rectangle. The general version only requires that the boundary of U be the image of the unit circle. Nevertheless, this version in this problem is pretty good.

13. Let U be a bounded open set for which Green's theorem holds. Let C be the oriented boundary consistent with Green's theorem. Show: area of $U = \int_C xdy$.
14. Let U be an open set with which is on one side of its boundary as above with the boundary being the image of a Lipschitz map which is one to one. (Note this condition eliminates the curve crossing itself.) Now suppose you have a closed polygonal curve going from $(x_0, y_0) \rightarrow (x_1, y_1) \rightarrow (x_2, y_2) \cdots (x_p, y_p) = (x_0, y_0)$ oriented such that the direction of motion is in the direction $\mathbf{k} \times \mathbf{n}$ where \mathbf{n} is the unit outer normal from the divergence theorem. Show that if U is this enclosed polygon, then

$$\text{Area of } U = \frac{1}{2} \sum_{k=1}^n (x_k + x_{k-1})(y_k - y_{k-1})$$

This is a pretty remarkable result. Just draw a few such polygons and think how you would find their area without it.

15. Orient the u, v axes just like the usual arrangement of the x and y axes, u axis like the x axis. Let $\mathbf{r} : U \rightarrow \mathbb{R}^3$ where U is an open subset of \mathbb{R}^2 . Suppose \mathbf{r} is C^2 and let \mathbf{F} be a C^1 vector field defined in V , an open set containing $\mathbf{r}(U)$. Show, using the above reduction identity of Problem 8 that

$$(\mathbf{r}_u \times \mathbf{r}_v) \cdot (\nabla \times \mathbf{F})(\mathbf{r}(u, v)) = ((\mathbf{F} \circ \mathbf{r})_u \cdot \mathbf{r}_v - (\mathbf{F} \circ \mathbf{r})_v \cdot \mathbf{r}_u)(u, v). \quad (18.15)$$

The left side is the dot product of the curl of the vector field \mathbf{F} with a normal vector to the surface $\mathbf{r}(U)$, namely $\mathbf{r}_u \times \mathbf{r}_v$. Show that the right side can be written as $((\mathbf{F} \circ \mathbf{r}) \cdot \mathbf{r}_v)_u - ((\mathbf{F} \circ \mathbf{r}) \cdot \mathbf{r}_u)_v$ thanks to equality of mixed partial derivatives. Now suppose U is a region for which Green's theorem holds, the curve C bounding U being Lipschitz. Verify Stokes' theorem

$$\int_U (\mathbf{r}_u \times \mathbf{r}_v) \cdot (\nabla \times \mathbf{F})(\mathbf{r}(u, v)) du dv = \int_{\mathbf{R}} (\mathbf{F} \circ \mathbf{r}) \cdot \mathbf{r}_u du + (\mathbf{F} \circ \mathbf{r}) \cdot \mathbf{r}_v dv$$

where $\mathbf{R}(t) = (u(t), v(t))$ for $t \rightarrow (u(t), v(t))$ being a parametrization of C oriented so that $(u'(t), v'(t), 0) \times \mathbf{k}$ is an exterior normal to U . Show the left integral can be written as $\int_{\mathbf{r}(U)} \nabla \times \mathbf{F} \cdot \mathbf{N} d\mathcal{H}^2$ where \mathbf{N} is a unit normal to the surface $\mathbf{r}(U)$. Show the integral on the right can be written as the differential form $\int_{\mathbf{R}} F_1 dx + F_2 dy + F_3 dz$ where $\mathbf{R}(t) = \mathbf{r}(u(t), v(t))$.

16. Let $f : \mathbb{C} \rightarrow \mathbb{C}$. Thus if $x + iy \in \mathbb{C}$, $f(x + iy) = u(x, y) + iv(x, y)$ where u is the real part and v is the imaginary part. As in one variable calculus, we define

$$\lim_{h \rightarrow 0} \frac{f(z+h) - f(z)}{h} = f'(z)$$

and we say that this derivative exists exactly when this limit exists. Consider $h = it$ and then let $h = t$. Take limits in these two ways and conclude that if $f'(z)$ exists, then it is given by

$$f'(z) = u_x + iv_x = v_y - iu_y$$

Thus you have the Cauchy Riemann equations $u_x = v_y, v_x = -u_y$. Show that if u, v are both C^1 , and these Cauchy Riemann equation hold, then the function will be differentiable. When this happens, we say the function is analytic. (In fact it can be shown that if the limit of the difference quotient exists, then these real and imaginary parts will automatically be continuous and the function will be analytic.)

17. For a function $f : \mathbb{C} \rightarrow \mathbb{C}$ which is continuous and $\gamma : [a, b] \rightarrow \Gamma \subseteq \mathbb{C}$ where Γ is a piecewise smooth curve, we define the contour integral

$$\int_{\Gamma} f(z) dz = \int_a^b f(\gamma(t)) \gamma'(t) dt.$$

Show that this equals $F(\gamma(b)) - F(\gamma(a))$ if f is analytic, this for some function F . In particular, if Γ is a suitable closed curve, then $\int_{\Gamma} f(z) dz = 0$. This is Cauchy's theorem from complex analysis.

18. Suppose $f \in L^1(U)$ where U is some open set in \mathbb{R}^p . Go ahead and assume f is Borel measurable although it should work with f only Lebesgue measurable. Show there is a set of m_{p-1} measure zero N such that if $\mathbf{x}_p \equiv (x_1, x_2, \dots, x_{p-1}) \notin N$, then $x_p \rightarrow f(\mathbf{x}_p, x_p)$ is in $L^1(U_{\mathbf{x}_p})$ where $U_{\mathbf{x}_p} = \{t : (\mathbf{x}_p, t) \in U\}$.
19. If $f \in L^1(U)$ and $f_{x_p} \in L^1(U)$ where $U \subseteq \mathbb{R}^p$ is a box like $\prod_k (a_k, b_k)$. Let f_{x_p} refer to the weak partial derivative. Can you show that for fixed $s, t \in (a_k, b_k)$ then for a.e. \mathbf{x}_p , $f(\mathbf{x}_p, t) - f(\mathbf{x}_p, s) = \int_s^t f_{x_p}(\mathbf{x}_p, \tau) d\tau$? Assume f is Borel measurable.

Part III

Abstract Theory

Chapter 19

Hausdorff Spaces and Measures

19.1 General Topological Spaces

It turns out that metric spaces are not sufficiently general for some applications. This section is a brief introduction to general topology. In making this generalization, the properties of balls in a metric space are stated as axioms for a subset of the power set of a given set X . This subset of the power set $\mathcal{P}(X)$ (set of all subsets) will be known as a basis for a topology. The properties of balls which are of importance are that the intersection of finitely many is the union of balls and that the union of all of them give the whole space. Recall that with a metric space, an open set was just one in which every point was an interior point. This simply meant that every point is contained in a ball which is contained in the given set. All that is being done here is to make these simple properties into axioms.

Definition 19.1.1 *Let X be a nonempty set and suppose $\mathcal{B} \subseteq \mathcal{P}(X)$. Then \mathcal{B} is a basis for a topology if it satisfies the following axioms.*

1.) *Whenever $p \in A \cap B$ for $A, B \in \mathcal{B}$, it follows there exists $C \in \mathcal{B}$ such that $p \in C \subseteq A \cap B$.*

2.) $\cup \mathcal{B} = X$.

Then a subset U , of X is an open set if for every point $x \in U$, there exists $B \in \mathcal{B}$ such that $x \in B \subseteq U$. Thus the open sets are exactly those which can be obtained as a union of sets of \mathcal{B} . Denote these subsets of X by the symbol τ and refer to τ as the topology or the set of open sets.

Note that this is simply the analog of saying a set is open exactly when every point is an interior point.

Proposition 19.1.2 *Let X be a set and let \mathcal{B} be a basis for a topology as defined above and let τ be the set of open sets determined by \mathcal{B} . Then*

$$\emptyset \in \tau, X \in \tau, \quad (19.1)$$

$$\text{If } \mathcal{C} \subseteq \tau, \text{ then } \cup \mathcal{C} \in \tau \quad (19.2)$$

$$\text{If } A, B \in \tau, \text{ then } A \cap B \in \tau. \quad (19.3)$$

Proof: If $p \in \emptyset$ then there exists $B \in \mathcal{B}$ such that $p \in B \subseteq \emptyset$ because there are no points in \emptyset . Therefore, $\emptyset \in \tau$. Now if $p \in X$, then by part 2.) of Definition 19.1.1 $p \in B \subseteq X$ for some $B \in \mathcal{B}$ and so $X \in \tau$.

If $\mathcal{C} \subseteq \tau$, and if $p \in \cup \mathcal{C}$, then there exists a set, $B \in \mathcal{C}$ such that $p \in B$. However, B is itself a union of sets from \mathcal{B} and so there exists $C \in \mathcal{B}$ such that $p \in C \subseteq B \subseteq \cup \mathcal{C}$. This verifies 19.2.

Finally, if $A, B \in \tau$ and $p \in A \cap B$, then since A and B are themselves unions of sets of \mathcal{B} , it follows there exists $A_1, B_1 \in \mathcal{B}$ such that $A_1 \subseteq A, B_1 \subseteq B$, and $p \in A_1 \cap B_1$. Therefore, by 1.) of Definition 19.1.1 there exists $C \in \mathcal{B}$ such that $p \in C \subseteq A_1 \cap B_1 \subseteq A \cap B$, showing that $A \cap B \in \tau$ as claimed. Of course from the above, if $A \cap B = \emptyset$, then $A \cap B \in \tau$. ■

Definition 19.1.3 *A set X together with such a collection of its subsets satisfying 19.1-19.3 is called a topological space. τ is called the topology or set of open sets of X .*

Definition 19.1.4 A topological space is said to be Hausdorff if whenever p and q are distinct points of X , there exist disjoint open sets U, V such that $p \in U$, $q \in V$. In other words points can be separated with open sets.



Definition 19.1.5 A subset of a topological space is said to be closed if its complement is open. Let p be a point of X and let $E \subseteq X$. Then p is said to be a limit point of E if every open set containing p contains a point of E distinct from p .

Theorem 19.1.6 If (X, τ) is a Hausdorff space and if $p \in X$, then $\{p\}$ is a closed set.

Proof: If $x \neq p$, there exist open sets U and V such that $x \in U$, $p \in V$ and $U \cap V = \emptyset$. Therefore, $\{p\}^C$ is an open set so $\{p\}$ is closed. ■

It would have been enough to assume that if $x \neq y$, then there exists an open set containing x which does not contain y .

Proposition 19.1.7 If (X, τ) is a Hausdorff space then a point p is a limit point of a set E if and only if every open set containing p contains infinitely many points of E each different than p .

Proof: \Leftarrow is obvious. Consider \Rightarrow . If p is a limit point and if U is an open set containing p but there are only finitely many points of E different than p contained in U , $\{q_i\}_{i=1}^m$, then consider $V \equiv U \cap \bigcap_{i=1}^m \{q_i\}^C$ which is an open set because each $\{q_i\}^C$ is open. This is because if $x \neq q_i$ there exists open V_{q_i} containing x such that $q_i \notin V_{q_i}$ and so V is a finite intersection of open sets. Therefore, there is a $q_{m+1} \in V \setminus \{p\}$, a contradiction. ■

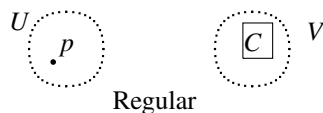
Theorem 19.1.8 A subset E , of X is closed if and only if it contains all its limit points. A set is closed if and only if its complement is open and a set is open if and only if its complement is closed.

Proof: Suppose first that E is closed and let x be a limit point of E . Is $x \in E$? If $x \notin E$, then E^C is an open set containing x which contains no points of E , a contradiction. Thus $x \in E$.

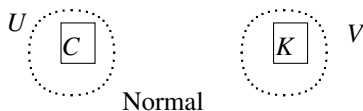
Now suppose E contains all its limit points. Is E^C open? If $x \in E^C$, then x is not a limit point of E because E has all its limit points and so there exists an open set, U containing x such that U contains no point of E other than x . Since $x \notin E$, it follows that $x \in U \subseteq E^C$ which implies E^C is an open set because this shows E^C is the union of open sets.

By definition, E closed $\Rightarrow E^C$ is open. If E^C is open, then no point of E^C can be a limit point of E and so E is closed since it contains all its limit points so E^C open $\Rightarrow E$ closed. ■

Definition 19.1.9 A topological space (X, τ) is said to be regular if whenever C is a closed set and p is a point not in C , there exist disjoint open sets U and V such that $p \in U$, $C \subseteq V$.



Definition 19.1.10 The topological space, (X, τ) is said to be normal if whenever C and K are disjoint closed sets, there exist disjoint open sets U, V with $C \subseteq U, K \subseteq V$. Thus any two disjoint closed sets can be separated with open sets.



Definition 19.1.11 Let E be a subset of X . \bar{E} is defined to be the smallest closed set containing E .

Lemma 19.1.12 The above definition is well defined.

Proof: Let \mathcal{C} denote all the closed sets which contain E . Then \mathcal{C} is nonempty because $X \in \mathcal{C}$.

$$(\cap \{A : A \in \mathcal{C}\})^C = \cup \{A^C : A \in \mathcal{C}\},$$

an open set which shows that $\cap \mathcal{C}$ is a closed set and is the smallest closed set which contains E . ■

Theorem 19.1.13 $\bar{E} = E \cup \{\text{limit points of } E\}$.

Proof: Let $x \in \bar{E}$ and suppose that $x \notin E$. If x is not a limit point either, then there exists an open set U , containing x which does not intersect E . But then U^C is a closed set which contains E which does not contain x , contrary to the definition that \bar{E} is the intersection of all closed sets containing E . Therefore, x must be a limit point of E after all.

Now $E \subseteq \bar{E}$ so suppose x is a limit point of E . Is $x \in \bar{E}$? If H is a closed set containing E , which does not contain x , then H^C is an open set containing x which contains no points of E other than x negating the assumption that x is a limit point of E . ■

The following is the definition of continuity in terms of general topological spaces. It is really just a generalization of the $\varepsilon - \delta$ definition of continuity given in calculus.

Definition 19.1.14 Let (X, τ) and (Y, η) be two topological spaces and let $f : X \rightarrow Y$. f is continuous at $x \in X$ if whenever V is an open set of Y containing $f(x)$, there exists an open set $U \in \tau$ such that $x \in U$ and $f(U) \subseteq V$. f is continuous if $f^{-1}(V) \in \tau$ whenever $V \in \eta$.

Then the following comes from the definition.

Proposition 19.1.15 In the situation of Definition 19.1.14 f is continuous if and only if f is continuous at every point of X .

Proof: \Rightarrow Suppose f is continuous and let $f(x) \in V$ an open set in Y . Then $x \in f^{-1}(V) \equiv U \in \tau$.

\Leftarrow Next suppose f is continuous at every point. Then if $V \in \eta$, and $x \in f^{-1}(V)$, continuity at x implies there is open $U_x \subseteq f^{-1}(V)$. Thus $f^{-1}(V) = \cup_{x \in f^{-1}(V)} U_x$ and so $f^{-1}(V)$ is open. ■

Definition 19.1.16 Let (X_i, τ_i) be topological spaces. $\prod_{i=1}^n X_i$ is the Cartesian product. Define a product topology as follows. Let $\mathcal{B} = \prod_{i=1}^n A_i$ where $A_i \in \tau_i$. Then \mathcal{B} is a basis for the product topology.

Theorem 19.1.17 The set \mathcal{B} of Definition 19.1.16 is a basis for a topology as claimed.

Proof: Suppose $x \in (\prod_{i=1}^n A_i) \cap (\prod_{i=1}^n B_i)$ where A_i and B_i are open sets. Suppose that $x = (x_1, \dots, x_n)$. Then $x_i \in A_i \cap B_i$ for each i . Therefore, $x \in \prod_{i=1}^n A_i \cap B_i \in \mathcal{B}$ and $\prod_{i=1}^n A_i \cap B_i \subseteq \prod_{i=1}^n A_i$. ■

The definition of compactness is also considered for a general topological space. This is given next.

Definition 19.1.18 A subset, E , of a topological space (X, τ) is said to be compact if whenever $\mathcal{C} \subseteq \tau$ and $E \subseteq \bigcup \mathcal{C}$, there exists a finite subset of \mathcal{C} , $\{U_1 \cdots U_n\}$, such that $E \subseteq \bigcup_{i=1}^n U_i$. (Every open covering admits a finite subcovering.) E is precompact if \bar{E} is compact. A topological space is called locally compact if it has a basis \mathcal{B} , with the property that \bar{B} is compact for each $B \in \mathcal{B}$.

In general topological spaces there may be no concept of “bounded”. Even if there is, closed and bounded is not necessarily the same as compactness. However, in any Hausdorff space every compact set must be a closed set.

Theorem 19.1.19 If (X, τ) is a Hausdorff space, then every compact subset must also be a closed set.

Proof: Suppose $p \notin K$ a compact set. For each $x \in X \setminus \{p\}$, there exist open sets, U_x and V_x such that $x \in U_x$, $p \in V_x$, and $U_x \cap V_x = \emptyset$. If K is assumed to be compact, there are finitely many of these sets, U_{x_1}, \dots, U_{x_m} which cover K . Then let $V \equiv \bigcap_{i=1}^m V_{x_i}$. It follows that V is an open set containing p which has empty intersection with each of the U_{x_i} . Consequently, V contains no points of K and is therefore not a limit point of K . ■

Definition 19.1.20 If every finite subset of a set P whose elements are sets has nonempty intersection, the set P is said to have the finite intersection property.

Theorem 19.1.21 Let \mathcal{K} be a set whose elements are compact subsets of a Hausdorff topological space, (X, τ) . Suppose \mathcal{K} has the finite intersection property. Then $\bigcap \mathcal{K} \neq \emptyset$.

Proof: Suppose to the contrary that $\emptyset = \bigcap \mathcal{K}$. Then consider $\mathcal{C} \equiv \{K^C : K \in \mathcal{K}\}$. It follows \mathcal{C} is an open cover of K_0 where K_0 is any particular element of \mathcal{K} . But then there are finitely many $K \in \mathcal{K}$, K_1, \dots, K_r such that $K_0 \subseteq \bigcup_{i=1}^r K_i^C$ implying that

$$K_0 \cap (\bigcap_{i=1}^r K_i) \subseteq (\bigcup_{i=1}^r K_i^C) \cap (\bigcap_{i=1}^r K_i) = (\bigcap_{i=1}^r K_i)^C \cap (\bigcap_{i=1}^r K_i) = \emptyset,$$

contradicting the finite intersection property. ■

There is a fundamental theorem, called Urysohn’s lemma which is valid for locally compact Hausdorff spaces which is presented next.

Lemma 19.1.22 Let X be a locally compact Hausdorff space and let $K \subseteq V \subseteq X$ where K is compact and V is open. Then there exists an open set U_k containing k such that \bar{U}_k is compact and $U_k \subseteq \bar{U}_k \subseteq V$. Also there exists U such that \bar{U} is compact and $K \subseteq U \subseteq \bar{U} \subseteq V$.

Proof: Since X is locally compact, there exists a basis of open sets whose closures are compact \mathcal{U} . Denote by \mathcal{C} the set of all $U \in \mathcal{U}$ which contain k and let \mathcal{C}' denote the set of all closures of these sets of \mathcal{C} intersected with the closed set V^C . Thus \mathcal{C}' is a collection of compact sets. There are finitely many of the sets of \mathcal{C}' which have empty intersection. If not, then \mathcal{C}' has the finite intersection property and so there exists a point p in all of them. Since X is a Hausdorff space, there exist disjoint basic open sets from \mathcal{U} , A, B such that $k \in A$ and $p \in B$. Therefore, $p \notin \bar{A}$ contrary to the above requirement that p be in all such sets. It follows there are sets A_1, \dots, A_m in \mathcal{C} such that $V^C \cap \bar{A}_1 \cap \dots \cap \bar{A}_m = \emptyset$. Let $U_k \equiv A_1 \cap \dots \cap A_m$. Then $\bar{U}_k \subseteq \bar{A}_1 \cap \dots \cap \bar{A}_m$ and so it has empty intersection with V^C . Thus it is contained in V . Also \bar{U}_k is a closed subset of the compact set \bar{A}_1 so it is compact and $k \in U_k$.

For the second part, consider all such U_k . Since K is compact, there are finitely many which cover K U_{k_1}, \dots, U_{k_n} . Then let $U \equiv \bigcup_{i=1}^n U_{k_i}$. It follows that $\bar{U} = \bigcup_{i=1}^n \bar{U}_{k_i}$ and each of these is compact so this set works. ■

The following is Urysohn's lemma for locally compact Hausdorff spaces.

Theorem 19.1.23 (Urysohn) *Let (X, τ) be locally compact and let $H \subseteq U$ where H is compact and U is open. Then there exists $g : X \rightarrow [0, 1]$ such that g is continuous, $g(x) = 1$ on H and $g(x) = 0$ if $x \notin V$ for some open set V such that $\bar{V} \subseteq U$ such that \bar{V} is compact.*

Proof: This involves using Lemma 19.1.22 repeatedly. First use this lemma to obtain V open such that its closure is compact and contained in U with $V \supseteq H$. Thus $H \subseteq V \subseteq \bar{V} \subseteq U, \bar{V}$ compact.

Let $D \equiv \{r_n\}_{n=1}^\infty$ be the rational numbers in $(0, 1)$. Using Lemma 19.1.22, let V_{r_1} be an open set such that

$$H \subseteq V_{r_1} \subseteq \bar{V}_{r_1} \subseteq V, \quad \bar{V}_{r_1} \text{ is compact}$$

Suppose V_{r_1}, \dots, V_{r_k} have been chosen and list the rational numbers r_1, \dots, r_k in order,

$$r_{l_1} < r_{l_2} < \dots < r_{l_k} \text{ for } \{l_1, \dots, l_k\} = \{1, \dots, k\}.$$

If $r_{k+1} > r_{l_k}$ then letting $p = r_{l_k}$, let $V_{r_{k+1}}$ satisfy

$$\bar{V}_p \subseteq V_{r_{k+1}} \subseteq \bar{V}_{r_{k+1}} \subseteq V, \quad \bar{V}_{r_{k+1}} \text{ compact}$$

If $r_{k+1} \in (r_{l_i}, r_{l_{i+1}})$, let $p = r_{l_i}$ and let $q = r_{l_{i+1}}$. Then let $V_{r_{k+1}}$ satisfy

$$\bar{V}_p \subseteq V_{r_{k+1}} \subseteq \bar{V}_{r_{k+1}} \subseteq V_q, \quad \bar{V}_{r_{k+1}} \text{ compact}$$

If $r_{k+1} < r_{l_1}$, let $p = r_{l_1}$ and let $V_{r_{k+1}}$ satisfy

$$H \subseteq V_{r_{k+1}} \subseteq \bar{V}_{r_{k+1}} \subseteq V_p, \quad \bar{V}_{r_{k+1}} \text{ compact}$$

Thus there exist open sets V_r for each $r \in \mathbb{Q} \cap (0, 1)$ with the property that if $r < s$,

$$H \subseteq V_r \subseteq \bar{V}_r \subseteq V_s \subseteq \bar{V}_s \subseteq V.$$

Now for $D \equiv \mathbb{Q} \cap (0, 1)$, in the following, t will be in D

$$f(x) \equiv \min(\inf\{t \in D : x \in V_t\}, 1), \quad f(x) \equiv 1 \text{ if } x \notin \bigcup_{t \in D} V_t.$$

I claim f is continuous. $f(x) < a$, means there must be some $t \in D$ with $t < a$ and $x \in V_t$. Thus $f^{-1}([0, a)) = \cup \{V_t : t < a, t \in D\}$, an open set.

Next consider $x \in f^{-1}([0, a])$ so $f(x) \leq a$. If $t > a$, then $x \in V_t$ from the definition of $f(x)$. Thus

$$f^{-1}([0, a]) \subseteq \cap \{V_t : t > a\} = \cap \{\bar{V}_t : t > a\}$$

which is a closed set. If $x \in \cap \{V_t : t > a\}$, then $f(x) \leq a$ from the definition and so equality holds and $f^{-1}([0, a])$ is closed. This is also true if $a = 1$. In this case, $f^{-1}([0, 1]) = X$. Therefore, for $a \geq 0$,

$$f^{-1}((a, 1]) \cup f^{-1}([0, a]) = X$$

and so $f^{-1}((a, 1])$ is an open set. It follows that f is continuous because $f^{-1}(a, b) = f^{-1}([0, b)) \cap f^{-1}((a, 1])$ the intersection of two open sets. Since this is so for every interval, it follows that the inverse image of any open set is open and so f is continuous. Clearly $f(x) = 0$ on H . If $x \in V^C$, then $x \notin V_t$ for any $t \in D$ so $f(x) = 1$ on V^C . Let $g(x) = 1 - f(x)$. ■

In any metric space there is a much easier proof of the conclusion of Urysohn's lemma which applies. The following is Lemma 3.12.1 listed here for convenience.

Lemma 19.1.24 *Let S be a nonempty subset of a metric space, (X, d) . Define*

$$f(x) \equiv \text{dist}(x, S) \equiv \inf \{d(x, y) : y \in S\}.$$

Then f is continuous.

In a metric space it is all much easier.

Theorem 19.1.25 *Let (X, τ) be a locally compact metric space in which the closures of balls are compact and let $H \subseteq U$ where H is compact and U is open. Then there exists $g : X \rightarrow [0, 1]$ such that g is continuous, $g(x) = 1$ on H and $g(x) = 0$ if $x \notin V$ for some open set V whose closure is compact.*

Proof: Let $\delta > 0$ be such that for all $h \in H, \text{dist}(h, U^C) > \delta$. This exists because $h \rightarrow \text{dist}(h, U^C)$ is continuous and so achieves its minimum on H which must be positive because U^C is closed. Now consider the balls $B(h, \delta)$. These cover the compact set H and so there are finitely many which do so. $B(h_1, \delta), \dots, B(h_m, \delta)$ where the closure of each of these is compact. Also $B(h_j, \delta) \subseteq U$. Because if $x \in \bar{B}(h_j, \delta)$, then $d(x, h_j) \leq \delta < \text{dist}(h_j, U^C)$. Let $V = \cup_{j=1}^m B(h_j, \delta)$. Thus $\bar{V} = \cup_{j=1}^m \bar{B}(h_j, \delta)$ because there are only finitely many sets. Also \bar{V} is compact because it is a finite union of compact sets. Now define

$$g(x) \equiv \frac{\text{dist}(x, V^C)}{\text{dist}(x, H) + \text{dist}(x, V^C)}$$

This is continuous, equals 1 on H and equals 0 off V because the denominator is always positive since both H, V^C are closed. ■

A useful construction when dealing with locally compact Hausdorff spaces is the notion of the one point compactification of the space.

Definition 19.1.26 *Suppose (X, τ) is a locally compact Hausdorff space. Then let $\tilde{X} \equiv X \cup \{\infty\}$ where ∞ is just the name of some point which is not in X which is called the point at infinity. A basis for the topology $\tilde{\tau}$ for \tilde{X} is*

$$\tau \cup \{K^C \text{ where } K \text{ is a compact subset of } X\}.$$

The complement is taken with respect to \tilde{X} and so the open sets, K^C are basic open sets which contain ∞ .

The reason this is called a compactification is contained in the next lemma.

Lemma 19.1.27 *If (X, τ) is a locally compact Hausdorff space, then $(\tilde{X}, \tilde{\tau})$ is a compact Hausdorff space. Also if U is an open set of $\tilde{\tau}$, then $U \setminus \{\infty\}$ is an open set of τ .*

Proof: Since (X, τ) is a locally compact Hausdorff space, it follows $(\tilde{X}, \tilde{\tau})$ is a Hausdorff topological space. The only case which needs checking is the one of $p \in X$ and ∞ . Since (X, τ) is locally compact, there exists an open set of τ , U having compact closure which contains p . Then $p \in U$ and $\infty \in \overline{U}^C$ and these are disjoint open sets containing the points, p and ∞ respectively. Now let \mathcal{C} be an open cover of \tilde{X} with sets from $\tilde{\tau}$. Then ∞ must be in some set, U_∞ from \mathcal{C} , which must contain a set of the form K^C where K is a compact subset of X . Then there exist sets from \mathcal{C} , U_1, \dots, U_r which cover K . Therefore, a finite subcover of \tilde{X} is $U_1, \dots, U_r, U_\infty$.

To see the last claim, suppose U contains ∞ since otherwise there is nothing to show. Notice that if C is a compact set, then $X \setminus C$ is an open set. Therefore, if $x \in U \setminus \{\infty\}$, and if $\tilde{X} \setminus C$ is a basic open set contained in U containing ∞ , then if x is in this basic open set of \tilde{X} , it is also in the open set $X \setminus C \subseteq U \setminus \{\infty\}$. If x is not in any basic open set of the form $\tilde{X} \setminus C$ then x is contained in an open set of τ which is contained in $U \setminus \{\infty\}$. Thus $U \setminus \{\infty\}$ is indeed open in τ . ■

Lemma 19.1.28 *Let (X, τ) be a topological space and let \mathcal{B} be a basis for τ . Then K is compact if and only if every open cover of basic open sets admits a finite subcover.*

Proof: Suppose first that X is compact. Then if \mathcal{C} is an open cover consisting of basic open sets, it follows it admits a finite subcover because these are open sets in \mathcal{C} .

Next suppose that every basic open cover admits a finite subcover and let \mathcal{C} be an open cover of X . Then define $\tilde{\mathcal{C}}$ to be the collection of basic open sets which are contained in some set of \mathcal{C} . It follows $\tilde{\mathcal{C}}$ is a basic open cover of X and so it admits a finite subcover, $\{U_1, \dots, U_p\}$. Now each U_i is contained in an open set of \mathcal{C} . Let O_i be a set of \mathcal{C} which contains U_i . Then $\{O_1, \dots, O_p\}$ is an open cover of X . ■

Actually, there is a profound generalization of this lemma.

19.2 The Alexander Sub-basis Theorem

The Hausdorff maximal theorem is one of several convenient versions of the axiom of choice. For a discussion of this, see the appendix on the subject. There is this one, the well ordering principal, and Zorn's lemma. They are all equivalent to the axiom of choice and which one you use is a matter of taste.

Theorem 19.2.1 *(Hausdorff maximal principle) Let \mathcal{F} be a nonempty partially ordered set. Then there exists a maximal chain.*

The main tool in the study of products of compact topological spaces is the Alexander subbasis theorem which is presented next. Recall a set is compact if every basic open cover admits a finite subcover, Lemma 19.1.28. This was pretty easy to prove. However, there is a much smaller set of open sets called a subbasis which has this property. The proof of this result is much harder.

Definition 19.2.2 $\mathcal{S} \subseteq \tau$ is called a subbasis for the topology τ if the set \mathcal{B} of finite intersections of sets of \mathcal{S} is a basis for the topology τ .

Theorem 19.2.3 Let (X, τ) be a topological space and let $\mathcal{S} \subseteq \tau$ be a subbasis for τ . Then if $H \subseteq X$, H is compact if and only if every open cover of H consisting entirely of sets of \mathcal{S} admits a finite subcover.

Proof: The only if part is obvious because the subbasic sets are themselves open.

If every basic open cover admits a finite subcover then the set in question is compact. Suppose then that H is a subset of X having the property that subbasic open covers admit finite subcovers. Is H compact? Assume this is not so. Then what was just observed about basic covers implies there exists a basic open cover of H , \mathcal{O} , which admits no finite subcover. Let \mathcal{F} be defined as

$$\{\mathcal{O} : \mathcal{O} \text{ is a basic open cover of } H \text{ which admits no finite subcover}\}.$$

The assumption is that \mathcal{F} is nonempty. Partially order \mathcal{F} by set inclusion and use the Hausdorff maximal principle to obtain a maximal chain, \mathcal{C} , of such open covers and let $\mathcal{D} = \cup \mathcal{C}$. If \mathcal{D} admits a finite subcover, then since \mathcal{C} is a chain and the finite subcover has only finitely many sets, some element of \mathcal{C} would also admit a finite subcover, contrary to the definition of \mathcal{F} . Therefore, \mathcal{D} admits no finite subcover. If \mathcal{D}' properly contains \mathcal{D} and \mathcal{D}' is a basic open cover of H , then \mathcal{D}' has a finite subcover of H since otherwise, \mathcal{C} would fail to be a maximal chain, being properly contained in $\mathcal{C} \cup \{\mathcal{D}'\}$. Every set of \mathcal{D} is of the form

$$U = \cap_{i=1}^m B_i, B_i \in \mathcal{S}$$

because they are all basic open sets. If it is the case that for all $U \in \mathcal{D}$ one of the B_i is found in \mathcal{D} , then replace each such U with the subbasic set from \mathcal{D} containing it. But then this would be a subbasic open cover of H which by assumption would admit a finite subcover contrary to the properties of \mathcal{D} . Therefore, one of the sets of \mathcal{D} , denoted by U , has the property that $U = \cap_{i=1}^m B_i$, $B_i \in \mathcal{S}$ and no B_i is in \mathcal{D} . Thus $\mathcal{D} \cup \{B_i\}$ admits a finite subcover, for each of the above B_i because it is strictly larger than \mathcal{D} . Let this finite subcover corresponding to B_i be denoted by $V_1^i, \dots, V_{m_i}^i, B_i$. Consider $\{U, V_j^i, j = 1, \dots, m_i, i = 1, \dots, m\}$. If $p \in H \setminus \cup \{V_j^i\}$, then $p \in B_i$ for each i and so $p \in U$. This is therefore a finite subcover of \mathcal{D} contradicting the properties of \mathcal{D} . Therefore, \mathcal{F} must be empty. ■

19.3 The Product Topology and Compactness

Now here is the definition of the product topological space.

Definition 19.3.1 Let (X_i, τ_i) for $i \in I$ be a topological space. By the axiom of choice, $\prod_{i \in I} X_i \neq \emptyset$. Then by definition, a basis for a topology on $\prod_{i \in I} X_i$ consists of sets of the form $\prod_{i \in I} A_i$ where $A_i = X_i$ except for **finitely** many i and for these, $A_i \in \tau_i$. A sub-basis for this topology consists of sets of the form $\prod_{i \in I} A_i$ where $A_i = X_i$ for all but a single i and for this one, $A_i \in \tau_i$. This product topology is denoted by $\prod \tau_i$. The resulting topological space is $(\prod_{i \in I} X_i, \prod \tau_i)$. “Subbasic” sets will be those which are in the sub-basis.

It is important that the basic open sets have the i^{th} entries not all of X_i in only **finitely many** i . If you don’t insist on this, you will be dealing with something called the “box

topology” and the following major theorem will not be true. You might go through the proof and see that this is the case. By the Alexander subbasis theorem, compactness is equivalent to saying that every open cover of subbasic sets admits a finite subcover.

Theorem 19.3.2 (Tychanoff) *If (X_i, τ_i) is compact, then so is $(\prod_{i \in I} X_i, \prod \tau_i)$.*

Proof: By the Alexander subbasis theorem, the theorem will be proved if every subbasic open cover admits a finite subcover. Therefore, let \mathcal{O} be a subbasic open cover of $\prod_{i \in I} X_i$. Let

$$\mathcal{O}_j = \{Q \in \mathcal{O} : Q = P_j(A) \text{ for some } A \in \tau_j\}.$$

Thus \mathcal{O}_j consists of those sets of \mathcal{O} which have a possibly proper subset of X_i only in the slot $i = j$. Let

$$\pi_j \mathcal{O}_j = \{A : P_j(A) \in \mathcal{O}_j\}.$$

Thus $\pi_j \mathcal{O}_j$ picks out those proper open subsets of X_j which occur in \mathcal{O}_j .

If no $\pi_j \mathcal{O}_j$ covers X_j , then by the axiom of choice, there exists $f \in \prod_{i \in I} X_i \setminus \cup \pi_i \mathcal{O}_i$. Therefore, $f(j) \notin \cup \pi_j \mathcal{O}_j$ for each $j \in I$. Now f is a point of $\prod_{i \in I} X_i$ and so $f \in P_k(A) \in \mathcal{O}$ for some k . However, this is a contradiction as it was shown that $f(k)$ is not an element of A . (A is one of the sets whose union makes up $\cup \pi_k \mathcal{O}_k$.) This contradiction shows that for some j , $\pi_j \mathcal{O}_j$ covers X_j . Thus $X_j = \cup \pi_j \mathcal{O}_j$ and so by compactness of X_j , there exist A_1, \dots, A_m , sets in τ_j such that $X_j \subseteq \cup_{i=1}^m A_i$ and $P_j(A_i) \in \mathcal{O}$. Therefore, $\{P_j(A_i)\}_{i=1}^m$ covers $\prod_{i \in I} X_i$. By the Alexander subbasis theorem this proves $\prod_{i \in I} X_i$ is compact. ■

19.4 Stone Weierstrass Theorem

This theorem was presented earlier in the context of a real algebra of functions on a compact set. Here this is extended to the case where the functions are defined on a locally compact Hausdorff space and also extended to the case where the functions have values in \mathbb{C} .

19.4.1 The Case of Locally Compact Sets

Definition 19.4.1 *Let (X, τ) be a locally compact Hausdorff space. $C_0(X)$ denotes the space of real or complex valued continuous functions defined on X with the property that if $f \in C_0(X)$, then for each $\varepsilon > 0$ there exists a compact set K such that $|f(x)| < \varepsilon$ for all $x \notin K$. Define $\|f\|_\infty = \sup \{|f(x)| : x \in X\}$.*

This norm is well defined because $|f(x)| < 1$ for x not in some compact set K and $|f(x)|$ achieves its maximum on K .

Lemma 19.4.2 *For (X, τ) a locally compact Hausdorff space with the above norm, $C_0(X)$ is a complete space.*

Proof: Let $(\tilde{X}, \tilde{\tau})$ be the one point compactification described in Lemma 19.1.27.

$$D \equiv \left\{ f \in C(\tilde{X}) : f(\infty) = 0 \right\}.$$

Then D is a closed subspace of $C(\tilde{X})$. For $f \in C_0(X)$,

$$\tilde{f}(x) \equiv \begin{cases} f(x) & \text{if } x \in X \\ 0 & \text{if } x = \infty \end{cases}$$

and let $\theta : C_0(X) \rightarrow D$ be given by $\theta f = \tilde{f}$. Then θ is one to one and onto and also satisfies $\|f\|_\infty = \|\theta f\|_\infty$. Now D is complete because it is a closed subspace of a complete space and so $C_0(X)$ with $\|\cdot\|_\infty$ is also complete. ■

The above refers to functions which have values in \mathbb{C} but the same proof works for functions which have values in any complete normed linear space.

In the case where the functions in $C_0(X)$ all have real values, I will denote the resulting space by $C_0(X; \mathbb{R})$ with similar meanings in other cases.

With this lemma, the generalization of the Stone Weierstrass theorem to locally compact sets is as follows.

Theorem 19.4.3 *Let \mathcal{A} be an algebra of functions in $C_0(X; \mathbb{R})$ where (X, τ) is a locally compact Hausdorff space which separates the points and annihilates no point. Then \mathcal{A} is dense in $C_0(X; \mathbb{R})$.*

Proof: Let $(\tilde{X}, \tilde{\tau})$ be the one point compactification as described in Lemma 19.1.27. Let $\tilde{\mathcal{A}}$ denote all finite linear combinations, $\left\{ \sum_{i=1}^n c_i \tilde{f}_i + c_0 : f \in \mathcal{A}, c_i \in \mathbb{R} \right\}$ where for $f \in C_0(X; \mathbb{R})$,

$$\tilde{f}(x) \equiv \begin{cases} f(x) & \text{if } x \in X \\ 0 & \text{if } x = \infty \end{cases}.$$

Then $\tilde{\mathcal{A}}$ is obviously an algebra of functions in $C(\tilde{X}; \mathbb{R})$. It separates points because this is true of \mathcal{A} . Similarly, it annihilates no point because of the inclusion of c_0 an arbitrary element of \mathbb{R} in the definition above. Therefore from Theorem 5.10.5, $\tilde{\mathcal{A}}$ is dense in $C(\tilde{X}; \mathbb{R})$. Letting $f \in C_0(X; \mathbb{R})$, it follows $\tilde{f} \in C(\tilde{X}; \mathbb{R})$ so there exists a sequence $\{h_n\} \subseteq \tilde{\mathcal{A}}$ such that h_n converges uniformly to \tilde{f} . Now h_n is of the form $\sum_{i=1}^n c_i^n \tilde{f}_i + c_0^n$ and since $\tilde{f}(\infty) = 0$, you can take each $c_0^n = 0$ and so this has shown the existence of a sequence of functions in \mathcal{A} such that it converges uniformly to f . ■

19.4.2 The Case of Complex Valued Functions

What about the general case where $C_0(X)$ consists of complex valued functions and the field of scalars is \mathbb{C} rather than \mathbb{R} ? The following is the version of the Stone Weierstrass theorem which applies to this case. You have to assume that for $f \in \mathcal{A}$ it follows $\bar{f} \in \mathcal{A}$.

Lemma 19.4.4 *Let z be a complex number. Then $\operatorname{Re}(z) = \operatorname{Im}(i\bar{z})$, $\operatorname{Im}(z) = \operatorname{Re}(i\bar{z})$.*

Proof: The following computation comes from the definition of real and imaginary parts.

$$\begin{aligned} \operatorname{Re}(z) &= \frac{z + \bar{z}}{2} = \frac{iz + i\bar{z}}{2i} = \frac{i\bar{z} - \overline{(i\bar{z})}}{2i} = \operatorname{Im}(i\bar{z}) \\ \operatorname{Im}(z) &= \frac{z - \bar{z}}{2i} = \frac{i\bar{z} - iz}{2} = \frac{i\bar{z} + \overline{(i\bar{z})}}{2} = \operatorname{Re}(i\bar{z}) \quad \blacksquare \end{aligned}$$

Theorem 19.4.5 *Suppose \mathcal{A} is an algebra of functions in $C_0(X)$ for X a locally compact Hausdorff space which separates the points of X and annihilates no point of X , and has the property that if $f \in \mathcal{A}$, then $\bar{f} \in \mathcal{A}$. Then \mathcal{A} is dense in $C_0(X)$.*

Proof: Let $\operatorname{Re} \mathcal{A} \equiv \{\operatorname{Re} f : f \in \mathcal{A}\}$, $\operatorname{Im} \mathcal{A} \equiv \{\operatorname{Im} f : f \in \mathcal{A}\}$.

Claim 1: $\operatorname{Re} \mathcal{A} = \operatorname{Im} \mathcal{A}$

Proof of claim: A typical element of $\operatorname{Re} \mathcal{A}$ is $\operatorname{Re} f$ where $f \in \mathcal{A}$, then from Lemma 19.4.4, $\operatorname{Re}(f) = \operatorname{Im}(i\bar{f}) \in \operatorname{Im} \mathcal{A}$. Thus $\operatorname{Re} \mathcal{A} \subseteq \operatorname{Im} \mathcal{A}$. By assumption, $i\bar{f} \in \mathcal{A}$. The other direction works the same. Just use the other formula in Lemma 19.4.4.

Claim 2: Both $\operatorname{Re} \mathcal{A}$ and $\operatorname{Im} \mathcal{A}$ are real algebras.

Proof of claim: It is obvious these are both real vector spaces. Since these are equal, it suffices to consider $\operatorname{Re} \mathcal{A}$. It remains to show that $\operatorname{Re} \mathcal{A}$ is closed with respect to products.

$$\frac{f + \bar{f}}{2} \frac{g + \bar{g}}{2} = \frac{1}{4} [fg + f\bar{g} + \bar{f}g + \bar{f}\bar{g}] = \frac{1}{4} [2\operatorname{Re}(fg) + 2\operatorname{Re}(\bar{f}g)]$$

Now by assumption, $fg \in \mathcal{A}$ and so $\operatorname{Re}(fg) \in \operatorname{Re} \mathcal{A}$. Also $\operatorname{Re}(\bar{f}g) \in \operatorname{Re} \mathcal{A}$ because both \bar{f}, g are in \mathcal{A} and it is an algebra. Thus, the above is in $\operatorname{Re} \mathcal{A}$ because, as noted, this is a real vector space.

Claim 3: $\mathcal{A} = \operatorname{Re} \mathcal{A} + i\operatorname{Im} \mathcal{A}$

Proof of claim: If $f \in \mathcal{A}$, then $f = \frac{f + \bar{f}}{2} + i\frac{f - \bar{f}}{2i} \in \operatorname{Re} \mathcal{A} + i\operatorname{Im} \mathcal{A}$ so $\mathcal{A} \subseteq \operatorname{Re} \mathcal{A} + i\operatorname{Im} \mathcal{A}$. Now a generic element of $\operatorname{Re} \mathcal{A} + i\operatorname{Im} \mathcal{A}$ is $\operatorname{Re}(f) + i\operatorname{Im}(g)$ for $f, g \in \mathcal{A}$.

$$\operatorname{Re}(f) + i\operatorname{Im}(g) \equiv \frac{f + \bar{f}}{2} + i\left(\frac{g - \bar{g}}{2i}\right) = \frac{f + g}{2} + \frac{\bar{f} - \bar{g}}{2} \in \mathcal{A}$$

because \mathcal{A} is closed with respect to conjugates. Thus $\operatorname{Re} \mathcal{A} + i\operatorname{Im} \mathcal{A} \subseteq \mathcal{A}$.

Both $\operatorname{Re} \mathcal{A}$ and $\operatorname{Im} \mathcal{A}$ must separate the points. Here is why: If $x_1 \neq x_2$, then there exists $f \in \mathcal{A}$ such that $f(x_1) \neq f(x_2)$. If $\operatorname{Im} f(x_1) \neq \operatorname{Im} f(x_2)$, this shows there is a function in $\operatorname{Im} \mathcal{A}$, $\operatorname{Im} f$ which separates these two points. If $\operatorname{Im} f$ fails to separate the two points, then $\operatorname{Re} f$ must separate the points and so, by Lemma 19.4.4,

$$\operatorname{Re} f(x_1) = \operatorname{Im}(i\bar{f}(x_1)) \neq \operatorname{Re} f(x_2) = \operatorname{Im}(i\bar{f}(x_2))$$

Thus $\operatorname{Im} \mathcal{A}$ separates the points. Similarly $\operatorname{Re} \mathcal{A}$ separates the points using a similar argument or because it is equal to $\operatorname{Im} \mathcal{A}$.

Neither $\operatorname{Re} \mathcal{A}$ nor $\operatorname{Im} \mathcal{A}$ annihilate any point. This is easy to see because if x is a point, there exists $f \in \mathcal{A}$ such that $f(x) \neq 0$. Thus either $\operatorname{Re} f(x) \neq 0$ or $\operatorname{Im} f(x) \neq 0$. If $\operatorname{Im} f(x) \neq 0$, this shows this point is not annihilated by $\operatorname{Im} \mathcal{A}$. Since they are equal, $\operatorname{Re} \mathcal{A}$ does not annihilate this point either.

It follows from Theorem 19.4.3 that $\operatorname{Re} \mathcal{A}$ and $\operatorname{Im} \mathcal{A}$ are dense in the real valued functions of $C_0(X)$. Let $f \in C_0(X)$. Then there exists $\{h_n\} \subseteq \operatorname{Re} \mathcal{A}$ and $\{g_n\} \subseteq \operatorname{Im} \mathcal{A}$ such that $h_n \rightarrow \operatorname{Re} f$ uniformly and $g_n \rightarrow \operatorname{Im} f$ uniformly. Therefore, $h_n + ig_n \in \mathcal{A}$ and it converges to f uniformly. ■

19.5 Partitions of Unity

As before, the idea of a partition of unity is of fundamental significance. It will be used to construct measures.

Definition 19.5.1 Define $\operatorname{spt}(f)$ (support of f) to be the closure of the set $\{x : f(x) \neq 0\}$. If V is an open set, $C_c(V)$ will be the set of continuous functions f , defined on Ω having $\operatorname{spt}(f) \subseteq V$. Thus in Theorem 19.1.23, $f \in C_c(V)$.

Definition 19.5.2 If K is a compact subset of an open set V , then $K \prec \phi \prec V$ if

$$\phi \in C_c(V), \phi(K) = \{1\}, \phi(\Omega) \subseteq [0, 1],$$

where Ω denotes the whole topological space considered. Also for $\phi \in C_c(\Omega)$, $K \prec \phi$ if

$$\phi(\Omega) \subseteq [0, 1] \text{ and } \phi(K) = 1.$$

and $\phi \prec V$ if $\phi \in C_c(V)$ and

$$\phi(\Omega) \subseteq [0, 1] \text{ and } \text{spt}(\phi) \subseteq V.$$

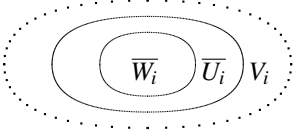
Theorem 19.5.3 (Partition of unity) Let K be a compact subset of a locally compact Hausdorff topological space satisfying Theorem 19.1.23 and suppose

$$K \subseteq V = \bigcup_{i=1}^n V_i, V_i \text{ open.}$$

Then there exist $\psi_i \prec V_i$ with $\sum_{i=1}^n \psi_i(x) = 1$ for all $x \in K$.

Proof: The proof is just like the one in Theorem 3.12.5 on Page 92. Let $K_1 = K \setminus \bigcup_{i=2}^n V_i$. Thus K_1 is compact and $K_1 \subseteq V_1$. Let $K_1 \subseteq W_1 \subseteq \overline{W}_1 \subseteq V_1$ with \overline{W}_1 compact. To obtain W_1 , use Theorem 19.1.23 to get f such that $K_1 \prec f \prec V_1$ and let $W_1 \equiv \{x : f(x) \neq 0\}$. Thus W_1, V_2, \dots, V_n covers K and $\overline{W}_1 \subseteq V_1$. Let $K_2 = K \setminus (\bigcup_{i=3}^n V_i \cup W_1)$. Then K_2 is compact and $K_2 \subseteq V_2$. Let $K_2 \subseteq W_2 \subseteq \overline{W}_2 \subseteq V_2$, \overline{W}_2 compact. Continue this way finally obtaining W_1, \dots, W_n , $K \subseteq W_1 \cup \dots \cup W_n$, and $\overline{W}_i \subseteq V_i$, \overline{W}_i compact. Now let $\overline{W}_i \subseteq U_i \subseteq \overline{U}_i \subseteq V_i$, \overline{U}_i compact.

By Theorem 19.1.23, let $\overline{U}_i \prec \phi_i \prec V_i$, $\bigcup_{i=1}^n \overline{W}_i \prec \gamma \prec \bigcup_{i=1}^n U_i$. Define



$$\psi_i(x) = \begin{cases} \gamma(x)\phi_i(x)/\sum_{j=1}^n \phi_j(x) & \text{if } \sum_{j=1}^n \phi_j(x) \neq 0, \\ 0 & \text{if } \sum_{j=1}^n \phi_j(x) = 0. \end{cases}$$

If x is such that $\sum_{j=1}^n \phi_j(x) = 0$, then $x \notin \bigcup_{i=1}^n \overline{U}_i$. Consequently $\gamma(y) = 0$ for all y near x and so $\psi_i(y) = 0$ for all y near x . Hence ψ_i is continuous at such x . If $\sum_{j=1}^n \phi_j(x) \neq 0$, this situation persists near x and so ψ_i is continuous at such points. Therefore ψ_i is continuous. If $x \in K$, then $\gamma(x) = 1$ and so $\sum_{j=1}^n \psi_j(x) = 1$. Clearly $0 \leq \psi_i(x) \leq 1$ and $\text{spt}(\psi_j) \subseteq V_j$. ■

The following corollary won't be needed immediately but is quite useful.

Corollary 19.5.4 In the context of the above theorem, if H is a compact subset of V_i , there exists a partition of unity such that $\psi_i(x) = 1$ for all $x \in H$ in addition to the conclusion of Theorem 19.5.3.

Proof: Keep V_i the same but replace V_j with $\tilde{V}_j \equiv V_j \setminus H$. Now in the proof above, applied to this modified collection of open sets, if $j \neq i$, $\phi_j(x) = 0$ whenever $x \in H$. Therefore, $\psi_i(x) = 1$ on H . ■

19.6 Measures on Hausdorff Spaces

In the case of a Hausdorff topological space, the following lemma gives conditions under which the σ algebra of μ measurable sets for an outer measure μ contains the Borel sets.

In words, it assumes the outer measure is inner regular on open sets and outer regular on all sets. Also it assumes you can approximate the measure of an open set with a compact set and the measure of a compact set with an open set. Recall that the Borel sets are those sets in the smallest σ algebra that contains the open sets. The big result is the Riesz representation theorem for positive linear functionals and the following lemma is the technical part of the proof of this big theorem in addition to being interesting for its own sake. It holds in a Hausdorff space, not just one which is locally compact.

Lemma 19.6.1 *Let Ω be a Hausdorff space and suppose μ is an outer measure satisfying μ is finite on compact sets and the following conditions,*

1. $\mu(E) = \inf \{ \mu(V), V \supseteq E, V \text{ open} \}$ for all E . (Outer regularity.)
2. For every open set V , $\mu(V) = \sup \{ \mu(K) : K \subseteq V, K \text{ compact} \}$ (Inner regularity on open sets.)
3. If A, B are compact disjoint sets, then $\mu(A \cup B) = \mu(A) + \mu(B)$.

Then the following hold.

1. If $\varepsilon > 0$ and if K is compact, there exists V open such that $V \supseteq K$ and

$$\mu(V \setminus K) < \varepsilon$$

2. If $\varepsilon > 0$ and if V is open with $\mu(V) < \infty$, there exists a compact subset K of V such that

$$\mu(V \setminus K) < \varepsilon$$

3. The μ measurable sets \mathcal{S} defined as

$$\mathcal{S} \equiv \{ E \subseteq \Omega : \mu(S) = \mu(S \setminus E) + \mu(S \cap E) \text{ for all } S \}$$

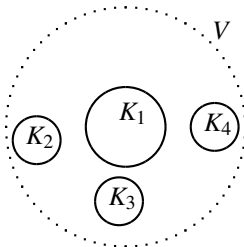
contains the Borel sets and also μ is inner regular on every open set,

$$\mu(V) = \sup \{ \mu(K) : K \subseteq V, K \text{ compact} \}$$

and for every $E \in \mathcal{S}$ with $\mu(E) < \infty$,

$$\mu(E) = \sup \{ \mu(K) : K \subseteq E, K \text{ compact} \}$$

Proof: First we establish 1 and 2 and use them to establish the last assertion. Consider 2. Suppose it is not true. Then there exists an open set V having $\mu(V) < \infty$ but for all $K \subseteq V$, $\mu(V \setminus K) \geq \varepsilon$ for some $\varepsilon > 0$. By inner regularity on open sets, there exists $K_1 \subseteq V$, K_1 compact, such that $\mu(K_1) \geq \varepsilon/2$. Now by assumption, $\mu(V \setminus K_1) \geq \varepsilon$ and so by inner regularity on open sets again, there exists compact $K_2 \subseteq V \setminus K_1$ such that $\mu(K_2) \geq \varepsilon/2$. Continuing this way, there is a sequence of disjoint compact sets contained in V $\{K_i\}$ such that $\mu(K_i) \geq \varepsilon/2$.



Now this is an obvious contradiction because by 3,

$$\mu(V) \geq \mu(\cup_{i=1}^n K_i) = \sum_{i=1}^n \mu(K_i) \geq n \frac{\varepsilon}{2}$$

for each n , contradicting $\mu(V) < \infty$.

Next consider 1. By outer regularity, there exists an open set $W \supseteq K$ such that $\mu(W) < \mu(K) + 1$. By 2, there exists compact $K_1 \subseteq W \setminus K$ such that $\mu((W \setminus K) \setminus K_1) < \varepsilon$. Then consider $V \equiv W \setminus K_1$. This is an open set containing K and from what was just shown,

$$\mu((W \setminus K_1) \setminus K) = \mu((W \setminus K) \setminus K_1) < \varepsilon.$$

Now consider the last assertion.

Define $\mathcal{S}_1 = \{E \in \mathcal{P}(\Omega) : E \cap K \in \mathcal{S}\}$ for all compact K .

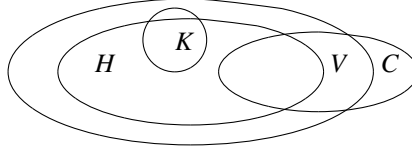
First it will be shown the compact sets are in \mathcal{S} . From this it will follow the closed sets are in \mathcal{S}_1 . Then you show $\mathcal{S}_1 = \mathcal{S}$. Thus $\mathcal{S}_1 = \mathcal{S}$ is a σ algebra and so it contains the Borel sets since it contains the closed sets. Finally you show the inner regularity assertion.

Claim 1: Compact sets are in \mathcal{S} .

Proof of claim: Let V be an open set with $\mu(V) < \infty$. I will show that for C compact,

$$\mu(V) \geq \mu(V \setminus C) + \mu(V \cap C).$$

If $\mu(V) = \infty$ the above is obvious. The various sets are illustrated in the following diagram.



By 2, there exists a compact set $K \subseteq V \setminus C$ such that $\mu((V \setminus C) \setminus K) < \varepsilon$ and a compact set $H \subseteq V$ such that $\mu(V \setminus H) < \varepsilon$. Thus $\mu(V) \leq \mu(V \setminus H) + \mu(H) < \varepsilon + \mu(H)$. Then

$$\begin{aligned} \mu(V) &\leq \mu(H) + \varepsilon \leq \mu(H \cap C) + \mu(H \setminus C) + \varepsilon \\ &\leq \mu(V \cap C) + \mu(V \setminus C) + \varepsilon \leq \mu(H \cap C) + \mu(K) + 3\varepsilon \end{aligned}$$

By 3,

$$= \mu(H \cap C) + \mu(K) + 3\varepsilon = \mu((H \cap C) \cup K) + 3\varepsilon \leq \mu(V) + 3\varepsilon.$$

Since ε is arbitrary, this shows that

$$\mu(V) = \mu(V \setminus C) + \mu(V \cap C). \quad (19.4)$$

Of course 19.4 is exactly what needs to be shown for arbitrary S in place of V . It suffices to consider only S having $\mu(S) < \infty$. If $S \subseteq \Omega$, with $\mu(S) < \infty$, let $V \supseteq S$, $\mu(S) + \varepsilon > \mu(V)$. Then from what was just shown, if C is compact,

$$\varepsilon + \mu(S) > \mu(V) = \mu(V \setminus C) + \mu(V \cap C) \geq \mu(S \setminus C) + \mu(S \cap C).$$

Since ε is arbitrary, this shows the compact sets are in \mathcal{S} . This proves the claim.

As discussed above, this verifies the closed sets are in \mathcal{S}_1 because if H is closed and C is compact, then compact $H \cap C \in \mathcal{S}$. If \mathcal{S}_1 is a σ algebra, this will show that \mathcal{S}_1 contains the Borel sets. Thus I first show \mathcal{S}_1 is a σ algebra.

To see that \mathcal{S}_1 is closed with respect to taking complements, let $E \in \mathcal{S}_1$ and K a compact set.

$$K = (E^C \cap K) \cup (E \cap K).$$

Then from the fact, just established, that the compact sets are in \mathcal{S} , $E^C \cap K = K \setminus (E \cap K) \in \mathcal{S}$. \mathcal{S}_1 is closed under countable unions because if K is a compact set and $E_n \in \mathcal{S}_1$,

$$K \cap \bigcup_{n=1}^{\infty} E_n = \bigcup_{n=1}^{\infty} K \cap E_n \in \mathcal{S}$$

because it is a countable union of sets of \mathcal{S} . Thus \mathcal{S}_1 is a σ algebra.

Therefore, if $E \in \mathcal{S}$ and K is a compact set, just shown to be in \mathcal{S} , it follows $K \cap E \in \mathcal{S}$ because \mathcal{S} is a σ algebra which contains the compact sets and so $\mathcal{S}_1 \supseteq \mathcal{S}$. It remains to verify $\mathcal{S}_1 \subseteq \mathcal{S}$. Recall that

$$\mathcal{S}_1 \equiv \{E : E \cap K \in \mathcal{S} \text{ for all } K \text{ compact}\}$$

Let $E \in \mathcal{S}_1$ and let V be an open set with $\mu(V) < \infty$ and choose $K \subseteq V$ such that $\mu(V \setminus K) < \varepsilon$. Then since $E \in \mathcal{S}_1$, it follows $E \cap K, E^C \cap K \in \mathcal{S}$ and so

$$\begin{aligned} \mu(V) &\leq \mu(V \setminus E) + \mu(V \cap E) \leq \overbrace{\mu(K \setminus E) + \mu(K \cap E)}^{\text{The two sets are disjoint and in } \mathcal{S}} + 2\varepsilon \\ &= \mu(K) + 2\varepsilon \leq \mu(V) + 3\varepsilon \end{aligned}$$

Since ε is arbitrary, this shows $\mu(V) = \mu(V \setminus E) + \mu(V \cap E)$ which would show $E \in \mathcal{S}$ if V were an arbitrary set.

Now let $S \subseteq \Omega$ be such an arbitrary set. If $\mu(S) = \infty$, then $\mu(S) = \mu(S \cap E) + \mu(S \setminus E)$. If $\mu(S) < \infty$, let

$$V \supseteq S, \mu(V) + \varepsilon \geq \mu(V).$$

Then

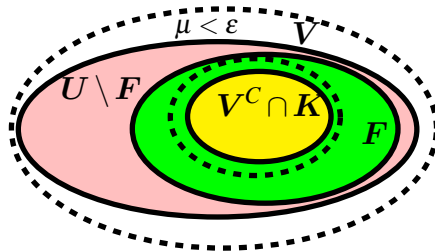
$$\mu(S) + \varepsilon \geq \mu(V) = \mu(V \setminus E) + \mu(V \cap E) \geq \mu(S \setminus E) + \mu(S \cap E).$$

Since ε is arbitrary, this shows that $E \in \mathcal{S}$ and so $\mathcal{S}_1 = \mathcal{S}$. Thus $\mathcal{S} \supseteq$ Borel sets as claimed.

From 2 μ is inner regular on all open sets. It remains to show that

$$\mu(F) = \sup\{\mu(K) : K \subseteq F\} \quad (19.5)$$

for all $F \in \mathcal{S}$ with $\mu(F) < \infty$. It might help to refer to the following crude picture to keep things straight. It also might not help. V is between the dotted lines.



From the picture as needed: Let $\mu(U \setminus F) < \varepsilon$ where U is open and let $K \subseteq U$ and $\mu(U \setminus K) < \varepsilon$, $\mu(V \setminus (U \setminus F)) < \varepsilon$ with V open and $V \supseteq U \setminus F = U \cap F^C$ so $V^C \subseteq U^C \cup F$. This is possible because all sets are in \mathcal{S} . Then $V^C \cap K \subseteq (U^C \cup F) \cap K = F \cap K \subseteq F$. Now $V^C \cap K$ is compact and

$$\begin{aligned} \mu(F \setminus (K \cap V^C)) &= \mu(F \cap (K^C \cup V)) = \mu(F \cap V) + \mu(F \cap K^C) \\ &\leq \mu(F \cap V) + \mu(U \setminus K) < \mu(F \cap V) + \varepsilon \end{aligned} \quad (19.6)$$

However,

$$\varepsilon > \mu(V \setminus (U \setminus F)) = \mu(V \cap (U \cap F^C)^C) = \mu(V \cap (U^C \cup F)) \geq \mu(V \cap F)$$

and so from 19.6, $\mu(F \setminus (K \cap V^C)) \leq 2\varepsilon$. Since $K \cap V^C$ is compact, this shows 19.5. ■

19.7 Measures and Positive Linear Functionals

This is on the Riesz representation theorem for positive linear functionals. It is a really marvelous result. It produces measures on locally compact Hausdorff spaces. Thus this doesn't help a lot in producing measures on infinite dimensional spaces but it works great on \mathbb{R}^n or closed subsets of \mathbb{R}^n and so forth.

Definition 19.7.1 Let (Ω, τ) be a topological space. $L : C_c(\Omega) \rightarrow \mathbb{C}$ is called a positive linear functional if L is linear, $L(af_1 + bf_2) = aLf_1 + bLf_2$, and if $Lf \geq 0$ whenever $f \geq 0$.

Theorem 19.7.2 (Riesz representation theorem) Let (Ω, τ) be a locally compact Hausdorff space and let L be a positive linear functional on $C_c(\Omega)$. Then there exists a σ algebra \mathcal{S} containing the Borel sets and a unique measure μ , defined on \mathcal{S} , such that

$$\mu \text{ is complete,} \quad (19.7)$$

$$\mu(K) < \infty \text{ for all } K \text{ compact,} \quad (19.8)$$

$$\mu(F) = \sup\{\mu(K) : K \subseteq F, K \text{ compact}\},$$

for all F open and for all $F \in \mathcal{S}$ with $\mu(F) < \infty$,

$$\mu(F) = \inf\{\mu(V) : V \supseteq F, V \text{ open}\}$$

for all $F \in \mathcal{S}$, and

$$\int f d\mu = Lf \text{ for all } f \in C_c(\Omega). \quad (19.9)$$

The plan is to define an outer measure and then to show that it, together with the σ algebra of sets measurable in the sense of Caratheodory, satisfies the conclusions of the theorem. Always, K will be a compact set and V will be an open set.

Definition 19.7.3 $\mu(V) \equiv \sup\{Lf : f \prec V\}$ for V open,

$$\mu(\emptyset) = 0, \mu(E) \equiv \inf\{\mu(V) : V \supseteq E\}$$

for arbitrary sets E .

Lemma 19.7.4 μ is a well-defined outer measure.

Proof: First it is necessary to verify that μ is well defined because there are two descriptions of it on open sets. Suppose then that $\mu_1(V) \equiv \inf\{\mu(U) : U \supseteq V \text{ and } U \text{ is open}\}$. It is required to verify that $\mu_1(V) = \mu(V)$ where μ is given as $\sup\{Lf : f \prec V\}$. If $U \supseteq V$, then $\mu(U) \geq \mu(V)$ directly from the definition. Hence from the definition of μ_1 , it follows $\mu_1(V) \geq \mu(V)$. On the other hand, $V \supseteq V$ and so $\mu_1(V) \leq \mu(V)$. This verifies μ is well defined.

It remains to show that μ is an outer measure. Let $V = \bigcup_{i=1}^{\infty} V_i$ and let $f \prec V$. Then $\text{spt}(f) \subseteq \bigcup_{i=1}^n V_i$ for some n . Let $\psi_i \prec V_i$, $\sum_{i=1}^n \psi_i = 1$ on $\text{spt}(f)$.

$$Lf = \sum_{i=1}^n L(f\psi_i) \leq \sum_{i=1}^n \mu(V_i) \leq \sum_{i=1}^{\infty} \mu(V_i).$$

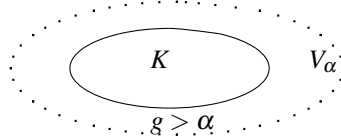
Hence $\mu(V) \leq \sum_{i=1}^{\infty} \mu(V_i)$ since $f \prec V$ is arbitrary. Now let $E = \bigcup_{i=1}^{\infty} E_i$. Is $\mu(E) \leq \sum_{i=1}^{\infty} \mu(E_i)$? Without loss of generality, it can be assumed $\mu(E_i) < \infty$ for each i since if not so, there is nothing to prove. Let $V_i \supseteq E_i$ with $\mu(E_i) + \varepsilon 2^{-i} > \mu(V_i)$.

$$\mu(E) \leq \mu(\bigcup_{i=1}^{\infty} V_i) \leq \sum_{i=1}^{\infty} \mu(V_i) \leq \varepsilon + \sum_{i=1}^{\infty} \mu(E_i).$$

Since ε was arbitrary, $\mu(E) \leq \sum_{i=1}^{\infty} \mu(E_i)$. It is clear from the definition that if $A \subseteq B$, then $\mu(A) \leq \mu(B)$. ■

Lemma 19.7.5 Let K be compact, $g \geq 0$, $g \in C_c(\Omega)$, and $g = 1$ on K . Then $\mu(K) \leq Lg$. Also $\mu(K) < \infty$ whenever K is compact.

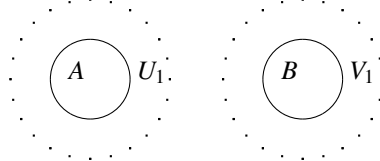
Proof: Let $\alpha \in (0, 1)$ and $V_\alpha = \{x : g(x) > \alpha\}$ so $V_\alpha \supseteq K$ and let $h \prec V_\alpha$.



Then $h \leq 1$ on V_α while $g\alpha^{-1} \geq 1$ on V_α and so $g\alpha^{-1} \geq h$ which implies $L(g\alpha^{-1}) \geq Lh$ and that therefore, since L is linear, $Lg \geq \alpha Lh$. Since $h \prec V_\alpha$ is arbitrary, and $K \subseteq V_\alpha$, $Lg \geq \alpha \mu(V_\alpha) \geq \alpha \mu(K)$. Letting $\alpha \uparrow 1$ yields $Lg \geq \mu(K)$. This proves the first part of the lemma. The second assertion follows from this and Theorem 19.1.23. If K is given, let $K \prec g \prec \Omega$ and so from what was just shown, $\mu(K) \leq Lg < \infty$. ■

Lemma 19.7.6 If A and B are disjoint compact subsets of Ω , then $\mu(A \cup B) = \mu(A) + \mu(B)$.

Proof: By Theorem 19.1.23, there exists $h \in C_c(\Omega)$ such that $A \prec h \prec B^C$. Let $U_1 = h^{-1}((\frac{1}{2}, 1])$, $V_1 = h^{-1}([0, \frac{1}{2}))$. Then $A \subseteq U_1$, $B \subseteq V_1$ and $U_1 \cap V_1 = \emptyset$.



From Lemma 19.7.5 $\mu(A \cup B) < \infty$ and so there exists an open set, W such that

$$W \supseteq A \cup B, \mu(A \cup B) + \varepsilon > \mu(W).$$

Now let $U = U_1 \cap W$ and $V = V_1 \cap W$. Then

$$U \supseteq A, V \supseteq B, U \cap V = \emptyset, \text{ and } \mu(A \cup B) + \varepsilon \geq \mu(W) \geq \mu(U \cup V).$$

Let $A \prec f \prec U, B \prec g \prec V$. Then by Lemma 19.7.5,

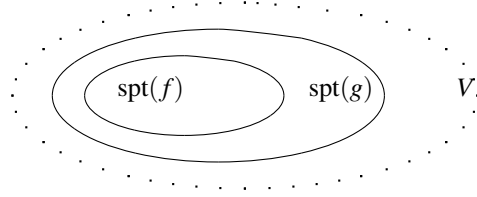
$$\mu(A \cup B) + \varepsilon \geq \mu(U \cup V) \geq L(f + g) = Lf + Lg \geq \mu(A) + \mu(B).$$

Since $\varepsilon > 0$ is arbitrary, this proves the lemma. ■

From Lemma 19.7.5 the following lemma is obtained.

Lemma 19.7.7 *Let $f \in C_c(\Omega)$, $f(\Omega) \subseteq [0, 1]$. Then $\mu(\text{spt}(f)) \geq Lf$. Also, every open set, V satisfies $\mu(V) = \sup \{\mu(K) : K \subseteq V\}$.*

Proof: Let $V \supseteq \text{spt}(f)$ and let $\text{spt}(f) \prec g \prec V$. Then $Lf \leq Lg \leq \mu(V)$ because $f \leq g$. Since this holds for all $V \supseteq \text{spt}(f)$, $Lf \leq \mu(\text{spt}(f))$ by definition of μ .



Finally, let V be open and let $l < \mu(V)$. Then from the definition of μ , there exists $f \prec V$ such that $L(f) > l$. Therefore, $l < \mu(\text{spt}(f)) \leq \mu(V)$ and so this shows the claim about inner regularity of the measure on an open set. ■

At this point, the conditions of Lemma 19.6.1 have been verified. Thus \mathcal{S} contains the Borel sets and μ is inner regular on sets of \mathcal{S} having finite measure.

It remains to show μ satisfies 19.9.

Lemma 19.7.8 $\int f d\mu = Lf$ for all $f \in C_c(\Omega)$.

Proof: Let $f \in C_c(\Omega)$, f real-valued, and suppose $f(\Omega) \subseteq [a, b]$. Choose $t_0 < a$ and let $t_0 < t_1 < \dots < t_n = b$, $t_i - t_{i-1} < \varepsilon$. Let

$$E_i = f^{-1}((t_{i-1}, t_i]) \cap \text{spt}(f). \quad (19.10)$$

Note that $\cup_{i=1}^n E_i = \text{spt}(f)$ since $\Omega = \cup_{i=1}^n f^{-1}((t_{i-1}, t_i])$. Let $V_i \supseteq E_i$, V_i is open and let V_i satisfy

$$f(x) < t_i + \varepsilon \text{ for all } x \in V_i, \mu(V_i \setminus E_i) < \varepsilon/n. \quad (19.11)$$

By Theorem 19.5.3 there exists $h_i \in C_c(\Omega)$ such that $h_i \prec V_i$, $\sum_{i=1}^n h_i(x) = 1$ on $\text{spt}(f)$. Now note that for each i , $f(x)h_i(x) \leq h_i(x)(t_i + \varepsilon)$. (If $x \in V_i$, this follows from 19.11. If $x \notin V_i$ both sides equal 0.) Therefore,

$$\begin{aligned} Lf &= L\left(\sum_{i=1}^n f h_i\right) \leq L\left(\sum_{i=1}^n h_i(t_i + \varepsilon)\right) = \sum_{i=1}^n (t_i + \varepsilon)L(h_i) \\ &= \sum_{i=1}^n (|t_0| + t_i + \varepsilon)L(h_i) - |t_0|L\left(\sum_{i=1}^n h_i\right). \end{aligned}$$

Now note that $|t_0| + t_i + \varepsilon \geq 0$ and so from the definition of μ and Lemma 19.7.5, this is no larger than

$$\begin{aligned}
 & \sum_{i=1}^n (|t_0| + t_i + \varepsilon) \mu(V_i) - |t_0| \mu(\text{spt}(f)) \\
 & \leq \sum_{i=1}^n (|t_0| + t_i + \varepsilon) (\mu(E_i) + \varepsilon/n) - |t_0| \mu(\text{spt}(f)) \\
 & \leq |t_0| \overbrace{\sum_{i=1}^n \mu(E_i)}^{\mu(\text{spt}(f))} + |t_0| \varepsilon + \sum_{i=1}^n t_i \mu(E_i) + \varepsilon(|t_0| + |b|) \\
 & \quad \sum_{i=1}^n t_i \frac{\varepsilon}{n} + \varepsilon \sum_{i=1}^n \mu(E_i) + \varepsilon^2 - |t_0| \mu(\text{spt}(f)).
 \end{aligned}$$

The first and last terms cancel. Therefore this is no larger than

$$\begin{aligned}
 & (2|t_0| + |b| + \mu(\text{spt}(f)) + \varepsilon) \varepsilon + \sum_{i=1}^n t_{i-1} \mu(E_i) + \varepsilon \mu(\text{spt}(f)) + \sum_{i=1}^n (|t_0| + |b|) \frac{\varepsilon}{n} \\
 & \leq \int f d\mu + (2|t_0| + |b| + 2\mu(\text{spt}(f)) + \varepsilon) \varepsilon + (|t_0| + |b|) \varepsilon
 \end{aligned}$$

Since $\varepsilon > 0$ is arbitrary, $Lf \leq \int f d\mu$ for all $f \in C_c(\Omega)$, f real valued. Hence equality holds because $L(-f) \leq -\int f d\mu$ so $L(f) \geq \int f d\mu$. Thus $Lf = \int f d\mu$ for all $f \in C_c(\Omega)$. Just apply the result for real functions to the real and imaginary parts of f . This gives the existence part of the Riesz representation theorem.

It only remains to prove uniqueness. Suppose both μ_1 and μ_2 are measures on \mathcal{S} satisfying the conclusions of the theorem. Then if K is compact and $V \supseteq K$, let $K \prec f \prec V$. Then

$$\mu_1(K) \leq \int f d\mu_1 = Lf = \int f d\mu_2 \leq \mu_2(V).$$

Thus, taking the inf for all $V \supseteq K$, $\mu_1(K) \leq \mu_2(K)$ for all K . Similarly, the inequality can be reversed and so it follows the two measures are equal on compact sets. By the assumption of inner regularity on open sets, the two measures are also equal on all open sets. By outer regularity, they are equal on all sets of \mathcal{S} . ■

Example 19.7.9 Let $L(f) = \int_{-\infty}^{\infty} f(t) dt$ for all $f \in C_c(\mathbb{R})$ where this is just the ordinary Riemann integral. Then the resulting measure is known as one dimensional Lebesgue measure.

Example 19.7.10 Let $L(f) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(x) dx_1 \cdots dx_n$ for $f \in C_c(\mathbb{R}^n)$. Then the resulting measure is m_n , n dimensional Lebesgue measure.

Here is a nice observation.

Proposition 19.7.11 In Example 19.7.10 the order of integration is not important. The same functional is obtained in any order.

Proof: Let $\text{spt}(f) \subseteq [-R, R]^n$. It clearly suffices to show this for $n = 2$. Then by the definition of the Riemann integral, $\int_{-R}^R \int_{-R}^R f(x, y) dx dy =$

$$\begin{aligned} \int_{-R}^R \sum_i \int_{x_i}^{x_{i+1}} f(s, y) ds dy &= \sum_{j=0}^{n-1} \sum_{i=0}^{n-1} \int_{y_i}^{y_{i+1}} \int_{x_i}^{x_{i+1}} f(s, t) ds dt \\ &= \sum_{j=0}^{n-1} \sum_{i=0}^{n-1} \int_{y_i}^{y_{i+1}} f(s_i, t) (x_{i+1} - x_i) dt \\ &= \sum_{j=0}^{n-1} \sum_{i=0}^{n-1} f(s_i, t_j) (x_{i+1} - x_i) (y_{j+1} - y_j) = \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} f(s_i, t_j) (x_{i+1} - x_i) (y_{j+1} - y_j) \end{aligned}$$

where $-R = x_0 < x_1 < \dots < x_n = R$ is a uniform partition of $[-R, R]$ with the y_i also giving a uniform partition of $[-R, R]$. Similar reasoning implies

$$\int_{-R}^R \int_{-R}^R f(x, y) dy dx = \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} f(\hat{s}_i, \hat{t}_j) (x_{i+1} - x_i) (y_{j+1} - y_j).$$

Now $(s_i, t_j), (\hat{s}_i, \hat{t}_j)$ are both in $[x_i, x_{i+1}] \times [y_j, y_{j+1}]$. Thus, by uniform continuity, if n is large enough,

$$|f(s_i, t_j) - f(\hat{s}_i, \hat{t}_j)| < \frac{\varepsilon}{4R^2}$$

Then it follows that $\left| \int_{-R}^R \int_{-R}^R f(x, y) dx dy - \int_{-R}^R \int_{-R}^R f(x, y) dy dx \right| \leq$

$$\begin{aligned} &\sum_{i=0}^{n-1} \sum_{j=0}^{n-1} |f(\hat{s}_i, \hat{t}_j) - f(s_i, t_j)| (x_{i+1} - x_i) (y_{j+1} - y_j) \\ &\leq \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} \frac{\varepsilon}{4R^2} (x_{i+1} - x_i) (y_{j+1} - y_j) = \varepsilon \end{aligned}$$

Since ε is arbitrary, this shows that the two iterated integrals are the same. In case $n > 2$, you can do exactly the same argument using the mean value theorem for integrals and obtain the same result by a similar argument, or you could use this result on pairs of integrals. ■

19.8 Slicing Measures

I saw this material first in the book [17]. It can be presented as an application of the theory of differentiation of Radon measures and the Riesz representation theorem for positive linear functionals. It is an amazing theorem and can be used to understand conditional probability. However, here I will obtain it from Theorem 10.14.12.

Theorem 19.8.1 *Let μ be a finite Radon measure on \mathbb{R}^{n+m} defined on a σ algebra, \mathcal{F} . Then there exists a unique finite Radon measure α , defined on a σ algebra \mathcal{S} , of sets of \mathbb{R}^n which satisfies*

$$\alpha(E) = \mu(E \times \mathbb{R}^m) \quad (19.12)$$

for all E Borel. There also exists a Borel set of α measure zero N , such that for each $x \notin N$, there exists a Radon probability measure ν_x such that if f is a nonnegative μ measurable function or a μ measurable function in $L^1(\mu)$,

$$y \rightarrow f(x, y) \text{ is } \nu_x \text{ measurable } \alpha \text{ a.e.}$$

$$x \rightarrow \int_{\mathbb{R}^m} f(x, y) d\nu_x(y) \text{ is } \alpha \text{ measurable} \quad (19.13)$$

and

$$\int_{\mathbb{R}^{n+m}} f(x, y) d\mu = \int_{\mathbb{R}^n} \left(\int_{\mathbb{R}^m} f(x, y) d\nu_x(y) \right) d\alpha(x). \quad (19.14)$$

If $\hat{\nu}_x$ is any other collection of Radon measures satisfying 19.13 and 19.14, then $\hat{\nu}_x = \nu_x$ for α a.e. x .

Proof: By Theorem 10.14.12 and the above lemmas, there exist unique Borel measurable α, ν_x such that 19.14 holds for all nonnegative Borel measurable functions f . This is because the Borel sets are contained in the product measurable sets. Now one can use the Riesz representation theorem on functionals $f \rightarrow \int_{\mathbb{R}^{n+m}} f(x, y) d\mu$ and $f \rightarrow \int_{\mathbb{R}^m} f(x, y) d\nu_x(y)$ along with regularity of these measures obtained from the Riesz representation theorem to extend and obtain the same result for f only μ measurable. ■

19.9 Exercises

1. Let X be a finite dimensional normed linear space, real or complex. Show that X is separable. **Hint:** Let $\{v_i\}_{i=1}^n$ be a basis and define a map from \mathbb{F}^n to X , θ , as follows. $\theta(\sum_{k=1}^n x_k e_k) \equiv \sum_{k=1}^n x_k v_k$. Show θ is continuous and has a continuous inverse. Now let D be a countable dense set in \mathbb{F}^n and consider $\theta(D)$.

2. Let $\alpha \in (0, 1]$. We define, for X a compact subset of \mathbb{R}^p ,

$$C^\alpha(X; \mathbb{R}^n) \equiv \{f \in C(X; \mathbb{R}^n) : \rho_\alpha(f) + \|f\| \equiv \|f\|_\alpha < \infty\}$$

where $\|f\| \equiv \sup\{|f(x)| : x \in X\}$ and

$$\rho_\alpha(f) \equiv \sup\left\{\frac{|f(x) - f(y)|}{|x - y|^\alpha} : x, y \in X, x \neq y\right\}.$$

Show that $(C^\alpha(X; \mathbb{R}^n), \|\cdot\|_\alpha)$ is a complete normed linear space. This is called a Holder space. What would this space consist of if $\alpha > 1$?

3. Let $\{f_n\}_{n=1}^\infty \subseteq C^\alpha(X; \mathbb{R}^n)$ where X is a compact subset of \mathbb{R}^p and suppose

$$\|f_n\|_\alpha \leq M$$

for all n . Show there exists a subsequence, n_k , such that f_{n_k} converges in $C(X; \mathbb{R}^n)$. We say the given sequence is precompact when this happens. (This also shows the embedding of $C^\alpha(X; \mathbb{R}^n)$ into $C(X; \mathbb{R}^n)$ is a compact embedding.) **Hint:** You might want to use the Ascoli Arzela theorem.

4. Suppose $f \in C_0([0, \infty))$ and also $|f(t)| \leq Ce^{-rt}$. Let \mathcal{A} denote the algebra of linear combinations of functions of the form e^{-st} for s sufficiently large. Thus \mathcal{A} is dense in $C_0([0, \infty))$. Show that if $\int_0^\infty e^{-st} f(t) dt = 0$ for each s sufficiently large, then $f(t) = 0$. Next consider only $|f(t)| \leq Ce^{rt}$ for some r . That is f has exponential growth. Show the same conclusion holds for f if $\int_0^\infty e^{-st} f(t) dt = 0$ for all s sufficiently large. This justifies the Laplace transform procedure of differential equations where if the Laplace transforms of two functions are equal, then the two functions are considered to be equal. More can be said about this. **Hint:** For the last part, consider $g(t) \equiv e^{-2rt} f(t)$ and apply the first part to g . If $g(t) = 0$ then so is $f(t)$.

5. A set S along with an order \leq is said to be a well ordered set if every nonempty subset of S has a smallest element. Here \leq is an order in the usual way. If $x, y \in S$, then either $x \leq y$ or $y \leq x$ and it satisfies the transitive law: If $x \leq y$ and $y \leq z$, then $x \leq z$. Using the Hausdorff maximal theorem, show that every nonempty set can be well ordered. That is, there is an order for which the given set is well ordered. In particular \mathbb{Q} the rational numbers can be well ordered. However, show that there is no way that the well order can coincide with the usual order on any open interval.
6. Verify that $P(\mathbb{N})$ the set of all subsets of the natural numbers is uncountable. Thus there exist uncountable sets. Pick an uncountable set Ω . Then you can consider this set to be well ordered, with the order denoted as \leq . Consider $\hat{\Omega} \equiv$

$$\{\omega \in \Omega \text{ such that there are uncountably many elements of } \Omega \text{ less than } \omega\}$$

If $\hat{\Omega} = \emptyset$, let $\Omega_0 = \hat{\Omega}$. If $\hat{\Omega} \neq \emptyset$, let ω_0 be the first element of $\hat{\Omega}$ and in this case let $\Omega_0 \equiv \{\omega : \omega < \omega_0\}$. That is $\omega \neq \omega_0$ and $\omega \leq \omega_0$. Explain why Ω_0 is uncountable. Explain why every element of Ω_0 has a “next” element. Now define a topology in the usual way. In particular, show that sets of the form $[\alpha, b), (a, b)$ where α is the first element of Ω_0 is a basis for a topology for Ω_0 . Verify that a sub-basis is sets of the form $[\alpha, b), (a, \omega_0)$. Explain why this is a Hausdorff space. Explain why every element of Ω_0 is preceded by countably many elements of Ω_0 and show every increasing sequence converges. Show that this cannot be a separable space. Suppose you have a cover of $(a, b]$ consisting of “sub-basic” open sets. Without loss of generality, all of these have nonempty intersection with (a, b) . Let p be the first such that (p, ω_0) is in the open cover, assuming there are such sets. Thus you could simply use (p, ω_0) instead of all the others. If $p \leq a$, you are done. If not, then $p \in (a, b]$ and so some set of the other kind, $[\alpha, q)$ must contain p and so at most two sets from the open cover contain (a, b) . Consider the other cases to verify that $(a, b]$ is compact. Now explain why (a, b) is either equal to $(a, b]$ or (a, b) . In the second case, verify that (a, b) would be of the form $(a, \hat{b}]$ which was just shown to be compact thanks to Alexander sub-basis theorem. Is Ω_0 locally compact?

7. In the above example, show that Ω_0 is not compact. However, show that every sequence has a subsequence which converges. Recall that in any metric space, compactness and sequential compactness are equivalent. Hence one can conclude that there is no metric which will deliver the same topology for Ω_0 described above. That is to say, this horrible topological space is not metrizable. However, the Riesz representation theorem presented above would hold for this terrible thing.

Chapter 20

Product Measures

Sometimes it is necessary to consider infinite Cartesian products of topological spaces.

20.1 Algebras

First of all, here is the definition of an algebra and theorems which tell how to recognize one when you see it. An algebra is like a σ algebra except it is only closed with respect to finite unions.

Definition 20.1.1 \mathcal{A} is said to be an algebra of subsets of a set Z if $Z \in \mathcal{A}$, $\emptyset \in \mathcal{A}$, and when $E, F \in \mathcal{A}$, $E \cup F$ and $E \setminus F$ are both in \mathcal{A} .

It is important to note that if \mathcal{A} is an algebra, then it is also closed under finite intersections. This is because $E \cap F = (E^C \cup F^C)^C \in \mathcal{A}$ since $E^C = Z \setminus E \in \mathcal{A}$ and $F^C = Z \setminus F \in \mathcal{A}$. Note that every σ algebra is an algebra but not the other way around.

Something satisfying the above definition is called an algebra because union is like addition, the set difference is like subtraction and intersection is like multiplication. Furthermore, only finitely many operations are done at a time and so there is nothing like a limit involved.

How can you recognize an algebra when you see one? The answer to this question is the purpose of the following lemma.

Lemma 20.1.2 Suppose \mathcal{R} and \mathcal{E} are subsets of $\mathcal{P}(Z)$ ¹ such that \mathcal{E} is defined as the set of all **finite disjoint unions** of sets of \mathcal{R} . Suppose also

$$\emptyset, Z \in \mathcal{R}$$

$$A \cap B \in \mathcal{R} \text{ whenever } A, B \in \mathcal{R},$$

$$A \setminus B \in \mathcal{E} \text{ whenever } A, B \in \mathcal{R}.$$

Then \mathcal{E} is an algebra of sets of Z .

Proof: Note first that if $A \in \mathcal{R}$, then $A^C \in \mathcal{E}$ because $A^C = Z \setminus A$. Now suppose that E_1 and E_2 are in \mathcal{E} ,

$$E_1 = \cup_{i=1}^m R_i, \quad E_2 = \cup_{j=1}^n R_j$$

where the R_i are disjoint sets in \mathcal{R} and the R_j are disjoint sets in \mathcal{R} . Then

$$E_1 \cap E_2 = \cup_{i=1}^m \cup_{j=1}^n R_i \cap R_j$$

which is clearly an element of \mathcal{E} because no two of the sets in the union can intersect and by assumption they are all in \mathcal{R} . Thus by induction, finite intersections of sets of \mathcal{E} are in \mathcal{E} . Consider the difference of two elements of \mathcal{E} next.

If $E = \cup_{i=1}^n R_i \in \mathcal{E}$, $E^C = \cap_{i=1}^n R_i^C =$ finite intersection of sets of \mathcal{E} which was just shown to be in \mathcal{E} . Now, if $E_1, E_2 \in \mathcal{E}$, $E_1 \setminus E_2 = E_1 \cap E_2^C \in \mathcal{E}$ from what was just shown about finite intersections.

¹Set of all subsets of Z

Finally consider finite unions of sets of \mathcal{E} . Let E_1 and E_2 be sets of \mathcal{E} . Then

$$E_1 \cup E_2 = (E_1 \setminus E_2) \cup E_2 \in \mathcal{E}$$

because $E_1 \setminus E_2$ consists of a finite disjoint union of sets of \mathcal{R} and these sets must be disjoint from the sets of \mathcal{R} whose union yields E_2 because $(E_1 \setminus E_2) \cap E_2 = \emptyset$. ■

The following corollary is particularly helpful in verifying the conditions of the above lemma.

Corollary 20.1.3 *Let $(Z_1, \mathcal{R}_1, \mathcal{E}_1)$ and $(Z_2, \mathcal{R}_2, \mathcal{E}_2)$ be as described in Lemma 20.1.2. Then $(Z_1 \times Z_2, \mathcal{R}, \mathcal{E})$ also satisfies the conditions of Lemma 20.1.2 if \mathcal{R} is defined as*

$$\mathcal{R} \equiv \{R_1 \times R_2 : R_i \in \mathcal{R}_i\}$$

and

$$\mathcal{E} \equiv \{\text{finite disjoint unions of sets of } \mathcal{R}\}.$$

Consequently, \mathcal{E} is an algebra of sets.

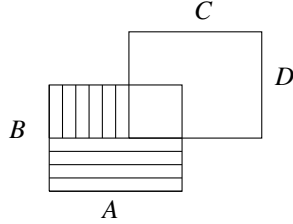
Proof: It is clear $\emptyset, Z_1 \times Z_2 \in \mathcal{R}$. Let $A \times B$ and $C \times D$ be two elements of \mathcal{R} .

$$A \times B \cap C \times D = (A \cap C) \times (B \cap D) \in \mathcal{R}$$

by assumption.

$$(A \times B) \setminus (C \times D) = A \times \overbrace{(B \setminus D)}^{\in \mathcal{E}_2} \cup \overbrace{(A \setminus C)}^{\in \mathcal{E}_1} \times \overbrace{(D \cap B)}^{\in \mathcal{R}_2} = (A \times Q) \cup (P \times R)$$

where $Q \in \mathcal{E}_2$, $P \in \mathcal{E}_1$, and $R \in \mathcal{R}_2$.



Since $A \times Q$ and $P \times R$ do not intersect, it follows that the above expression is in \mathcal{E} because each of these terms are. ■

20.2 Caratheodory Extension Theorem

The Caratheodory extension theorem is a fundamental result which makes possible the consideration of measures on infinite products among other things. The idea is that if a finite measure defined only on an algebra is trying to be a measure, then in fact it can be extended to a measure.

Definition 20.2.1 *Let \mathcal{E} be an algebra of sets of Ω . Thus \mathcal{E} contains \emptyset, Ω , and is closed with respect to differences and finite unions. Then μ_0 is a finite measure on \mathcal{E} means μ_0 is finitely additive: If E_i, E are sets of \mathcal{E} with the E_i disjoint and $E = \bigcup_{i=1}^{\infty} E_i$, then $\mu(E) = \sum_{i=1}^{\infty} \mu(E_i)$*

In this definition, μ_0 is trying to be a measure and acts like one whenever possible. Note the extra assumption that $E \in \mathcal{E}$. This would be automatic if it were a finite sum and finite union. Under these conditions, μ_0 can be extended uniquely to a complete measure μ , defined on a σ algebra of sets containing \mathcal{E} such that μ agrees with μ_0 on \mathcal{E} . The following is the main result.

Theorem 20.2.2 *Let μ_0 be a measure on an algebra of sets \mathcal{E} , which satisfies $\mu_0(\Omega) < \infty$. Then there exists a complete measure space $(\Omega, \mathcal{S}, \mu)$ such that $\mu(E) = \mu_0(E)$ for all $E \in \mathcal{E}$. Also if ν is any measure which agrees with μ_0 on \mathcal{E} , then $\nu = \mu$ on $\sigma(\mathcal{E})$, the σ algebra generated by \mathcal{E} .*

Proof: Define an outer measure as follows.

$$\mu(S) \equiv \inf \left\{ \sum_{i=1}^{\infty} \mu_0(E_i) : S \subseteq \bigcup_{i=1}^{\infty} E_i, E_i \in \mathcal{E} \right\}$$

Claim 1: μ is an outer measure.

Proof of Claim 1: Let $S \subseteq \bigcup_{i=1}^{\infty} S_i$ and let $S_i \subseteq \bigcup_{j=1}^{\infty} E_{ij}$, where

$$\mu(S_i) + \frac{\varepsilon}{2^i} \geq \sum_{j=1}^{\infty} \mu(E_{ij}).$$

Then

$$\mu(S) \leq \sum_i \sum_j \mu(E_{ij}) = \sum_i \left(\mu(S_i) + \frac{\varepsilon}{2^i} \right) = \sum_i \mu(S_i) + \varepsilon.$$

Since ε is arbitrary, this shows μ is an outer measure as claimed.

By the Caratheodory procedure, there exists a unique σ algebra \mathcal{S} , consisting of the μ measurable sets such that $(\Omega, \mathcal{S}, \mu)$ is a complete measure space. It remains to show that μ extends μ_0 .

Claim 2: If \mathcal{S} is the σ algebra of μ measurable sets, $\mathcal{S} \supseteq \mathcal{E}$ and $\mu = \mu_0$ on \mathcal{E} .

Proof of Claim 2: First observe that if $A \in \mathcal{E}$, then $\mu(A) \leq \mu_0(A)$ by definition. It remains to turn the inequality around. Letting

$$\mu(A) + \varepsilon > \sum_{i=1}^{\infty} \mu_0(E_i), \bigcup_{i=1}^{\infty} E_i \supseteq A, E_i \in \mathcal{E},$$

it follows that $\mu(A) + \varepsilon > \sum_{i=1}^{\infty} \mu_0(E_i \cap A) \geq \mu_0(A)$ since $A = \bigcup_{i=1}^{\infty} E_i \cap A$. Therefore, $\mu = \mu_0$ on \mathcal{E} .

Consider the assertion that $\mathcal{E} \subseteq \mathcal{S}$. Let $A \in \mathcal{E}$ and let $S \subseteq \Omega$ be any set. By definition, there exist sets $\{E_i\} \subseteq \mathcal{E}$ such that $\bigcup_{i=1}^{\infty} E_i \supseteq S$ but

$$\mu(S) + \varepsilon > \sum_{i=1}^{\infty} \mu(E_i) = \sum_{i=1}^{\infty} \mu_0(E_i).$$

Then by the assumption that μ_0 is a measure on \mathcal{E} and that $\mu = \mu_0$ on sets of \mathcal{E} ,

$$\begin{aligned} \mu(S) &\leq \mu(S \cap A) + \mu(S \setminus A) \\ &\leq \mu((\bigcup_{i=1}^{\infty} E_i) \setminus A) + \mu(\bigcup_{i=1}^{\infty} (E_i \cap A)) \leq \sum_{i=1}^{\infty} \mu(E_i \setminus A) + \sum_{i=1}^{\infty} \mu(E_i \cap A) \end{aligned}$$

$$= \sum_{i=1}^{\infty} \mu_0(E_i \setminus A) + \sum_{i=1}^{\infty} \mu_0(E_i \cap A) = \sum_{i=1}^{\infty} \mu(E_i) < \mu(S) + \varepsilon.$$

Since ε is arbitrary, this shows $A \in \mathcal{S}$.

This has proved the existence part of the theorem. To verify uniqueness, $\sigma(\mathcal{E}) \subseteq \mathcal{S}$. Let $\mathcal{G} \equiv \{E \in \sigma(\mathcal{E}) : \mu(E) = \nu(E)\}$. Then \mathcal{G} is given to contain \mathcal{E} and is obviously closed with respect to countable disjoint unions and complements because $\sigma(\mathcal{E}) \subseteq \mathcal{S}$. Therefore by Dynkin's theorem, Lemma 9.3.2, $\mathcal{G} = \sigma(\mathcal{E})$. ■

The following lemma is also very significant. Actually Lemmas 9.8.4 and 9.8.5 are of even more use, but one can use the following to get useful information in some cases. In particular, in the proof of the Kolmogorov extension theorem, one can consider the special case that $M_t = \mathbb{R}$ or \mathbb{R}^{n_t} in that theorem. The following lemma is Corollary 9.8.9 on Page 9.8.9. However, I am giving a different proof here to emphasize the theorem on positive linear functionals and measures.

Lemma 20.2.3 *Let M be a metric space with the closed balls compact and suppose μ is a measure defined on the Borel sets of M which is finite on compact sets. Then there exists a unique Radon measure, $\bar{\mu}$ which equals μ on the Borel sets. In particular μ must be both inner and outer regular on all Borel sets.*

Proof: Define a positive linear functional, $\Lambda(f) = \int f d\mu$. Let $\bar{\mu}$ be the Radon measure which comes from the Riesz representation theorem for positive linear functionals. Thus for all $f \in C_c(M)$, $\int f d\mu = \int f d\bar{\mu}$. If V is an open set, let $\{f_n\}$ be a sequence of continuous functions in $C_c(M)$ which is increasing and converges to \mathcal{X}_V pointwise. Then applying the monotone convergence theorem, $\int \mathcal{X}_V d\mu = \mu(V) = \int \mathcal{X}_V d\bar{\mu} = \bar{\mu}(V)$ and so the two measures coincide on all open sets. Every compact set is a countable intersection of open sets and so the two measures coincide on all compact sets. Now let $B(a, n)$ be a ball of radius n and let E be a Borel set contained in this ball. Then by regularity of $\bar{\mu}$ there exist sets F, G such that G is a countable intersection of open sets and F is a countable union of compact sets such that $F \subseteq E \subseteq G$ and $\bar{\mu}(G \setminus F) = 0$. Now $\mu(G) = \bar{\mu}(G)$ and $\mu(F) = \bar{\mu}(F)$. Thus

$$\bar{\mu}(G \setminus F) + \bar{\mu}(F) = \bar{\mu}(G) = \mu(G) = \mu(G \setminus F) + \mu(F)$$

and so $\mu(G \setminus F) = \bar{\mu}(G \setminus F)$. It follows $\mu(E) = \mu(F) = \bar{\mu}(F) = \bar{\mu}(G) = \bar{\mu}(E)$. If E is an arbitrary Borel set, then $\mu(E \cap B(a, n)) = \bar{\mu}(E \cap B(a, n))$ and letting $n \rightarrow \infty$, this yields $\mu(E) = \bar{\mu}(E)$. ■

20.3 Kolmogorov Extension Theorem

This extension theorem is one of the most important theorems in probability theory, at least according to my understanding of the situation. As an example, one sometimes wants to consider infinitely many independent normally distributed random variables. Is there a probability space such that this kind of thing even exists? The answer is yes and one way to show this is through the use of the Kolmogorov extension theorem. I am presenting the most general version of this theorem that I have seen. For another proof see the book by Strook [56]. What I am using here is a modification of one in Billingsley [6].

Let M_t be a complete separable metric space. This is called a Polish space. I will denote a totally ordered index set, (Like \mathbb{R}) and the interest will be in building a measure on the product space, $\prod_{t \in I} M_t$. If you like less generality, just think of $M_t = \mathbb{R}^{k_t}$ or even

$M_t = \mathbb{R}$. By the well ordering principle, you can always put an order on any index set so this order is no restriction, but we do not insist on a well order and in fact, index sets of most interest are \mathbb{R} or $[0, \infty)$. Also for X a topological space, $\mathcal{B}(X)$ will denote the Borel sets.

Notation 20.3.1 The symbol J will denote a finite subset of I , $J = (t_1, \dots, t_n)$, the t_i taken in order. \mathcal{E}_J will denote a set which has a set E_t of $\mathcal{B}(M_t)$ in the t^{th} position for $t \in J$ and for $t \notin J$, the set in the t^{th} position will be M_t . \mathcal{K}_J will denote a set which has a compact set in the t^{th} position for $t \in J$ and for $t \notin J$, the set in the t^{th} position will be M_t . Also denote by \mathcal{R}_J the sets \mathcal{E}_J and \mathcal{R} the union of all such \mathcal{R}_J . Let \mathcal{E}_J denote finite disjoint unions of sets of \mathcal{R}_J and let \mathcal{E} denote finite disjoint unions of sets of \mathcal{R} . Thus if \mathbf{F} is a set of \mathcal{E} , there exists J such that \mathbf{F} is a finite disjoint union of sets of \mathcal{R}_J . For $\mathbf{F} \in \Omega$, denote by $\pi_J(\mathbf{F})$ the set $\prod_{t \in J} F_t$ where $\mathbf{F} = \prod_{t \in I} F_t$.

Lemma 20.3.2 The sets $\mathcal{E}, \mathcal{E}_J$ defined above form an algebra of sets of $\prod_{t \in I} M_t$.

Proof: First consider \mathcal{R}_J . If $\mathbf{A}, \mathbf{B} \in \mathcal{R}_J$, then $\mathbf{A} \cap \mathbf{B} \in \mathcal{R}_J$ also. Is $\mathbf{A} \setminus \mathbf{B}$ a finite disjoint union of sets of \mathcal{R}_J ? It suffices to verify that $\pi_J(\mathbf{A} \setminus \mathbf{B})$ is a finite disjoint union of $\pi_J(\mathcal{R}_J)$. Let $|J|$ denote the number of indices in J . If $|J| = 1$, then it is obvious that $\pi_J(\mathbf{A} \setminus \mathbf{B})$ is a finite disjoint union of sets of $\pi_J(\mathcal{R}_J)$. In fact, letting $J = (t)$ and the t^{th} entry of \mathbf{A} is A and the t^{th} entry of \mathbf{B} is B , then the t^{th} entry of $\mathbf{A} \setminus \mathbf{B}$ is $A \setminus B$, a Borel set of M_t , a finite disjoint union of Borel sets of M_t .

Suppose then that for \mathbf{A}, \mathbf{B} sets of \mathcal{R}_J , $\pi_J(\mathbf{A} \setminus \mathbf{B})$ is a finite disjoint union of sets of $\pi_J(\mathcal{R}_J)$ for $|J| \leq n$, and consider $J = (t_1, \dots, t_n, t_{n+1})$. Let the t_i^{th} entry of \mathbf{A} and \mathbf{B} be respectively A_i and B_i . It follows that $\pi_J(\mathbf{A} \setminus \mathbf{B})$ has the following in the entries for J

$$(A_1 \times A_2 \times \dots \times A_n \times A_{n+1}) \setminus (B_1 \times B_2 \times \dots \times B_n \times B_{n+1})$$

Letting A represent $A_1 \times A_2 \times \dots \times A_n$ and B represent $B_1 \times B_2 \times \dots \times B_n$, $\mathbf{A} \setminus \mathbf{B}$ is of the form

$$A \times (A_{n+1} \setminus B_{n+1}) \cup (A \setminus B) \times (A_{n+1} \cap B_{n+1})$$

By induction, $(A \setminus B)$ is the finite disjoint union of sets of $\mathcal{R}_{(t_1, \dots, t_n)}$. Therefore, the above is the finite disjoint union of sets of \mathcal{R}_J . It follows that \mathcal{E}_J is an algebra.

Now suppose $\mathbf{A}, \mathbf{B} \in \mathcal{R}$. Then for some finite set J , both are in \mathcal{R}_J . Then from what was just shown,

$$\mathbf{A} \setminus \mathbf{B} \in \mathcal{E}_J \subseteq \mathcal{E}, \mathbf{A} \cap \mathbf{B} \in \mathcal{R}.$$

By Lemma 20.1.2 on Page 523 this shows \mathcal{E} is an algebra. ■

With this preparation, here is the Kolmogorov extension theorem. In the statement and proof of the theorem, F_i, G_i , and E_i will denote Borel sets. Any list of indices from I will always be assumed to be taken in order. Thus, if $J \subseteq I$ and $J = (t_1, \dots, t_n)$, it will always be assumed $t_1 < t_2 < \dots < t_n$.

Theorem 20.3.3 For each finite set

$$J = (t_1, \dots, t_n) \subseteq I,$$

suppose there exists a Borel probability measure, $\nu_J = \nu_{t_1, \dots, t_n}$ defined on the Borel sets of $\prod_{t \in J} M_t$ such that the following consistency condition holds. If

$$(t_1, \dots, t_n) \subseteq (s_1, \dots, s_p),$$

then

$$\nu_{t_1 \cdots t_n}(F_{t_1} \times \cdots \times F_{t_n}) = \nu_{s_1 \cdots s_p}(G_{s_1} \times \cdots \times G_{s_p}) \quad (20.1)$$

where if $s_i = t_j$, then $G_{s_i} = F_{t_j}$ and if s_i is not equal to any of the indices t_k , then $G_{s_i} = M_{s_i}$. Then for \mathcal{E} defined in Notation 20.3.1, there exists a probability measure P and a σ algebra $\mathcal{F} = \sigma(\mathcal{E})$ such that $(\prod_{t \in I} M_t, P, \mathcal{F})$ is a probability space. Also there exist measurable functions, $X_s : \prod_{t \in I} M_t \rightarrow M_s$ defined for $s \in I$ as $X_s \equiv x_s$ such that for each $(t_1 \cdots t_n) \subseteq I$,

$$\begin{aligned} \nu_{t_1 \cdots t_n}(F_{t_1} \times \cdots \times F_{t_n}) &= P([X_{t_1} \in F_{t_1}] \cap \cdots \cap [X_{t_n} \in F_{t_n}]) \\ &= P\left((X_{t_1}, \dots, X_{t_n}) \in \prod_{j=1}^n F_{t_j}\right) = P\left(\prod_{t \in I} F_t\right) \end{aligned} \quad (20.2)$$

where $F_t = M_t$ for every $t \notin \{t_1 \cdots t_n\}$ and F_{t_i} is a Borel set. Also if f is a nonnegative function of finitely many variables, x_{t_1}, \dots, x_{t_n} , measurable with respect to $\mathcal{B}(\prod_{j=1}^n M_{t_j})$, then f is also measurable with respect to \mathcal{F} and

$$\int_{M_{t_1} \times \cdots \times M_{t_n}} f(x_{t_1}, \dots, x_{t_n}) d\nu_{t_1 \cdots t_n} = \int_{\prod_{t \in I} M_t} f(x_{t_1}, \dots, x_{t_n}) dP \quad (20.3)$$

Proof: Let \mathcal{E} be the algebra of sets defined in the above notation. I want to define a measure on \mathcal{E} . For $F \in \mathcal{E}$, there exists J such that F is the finite disjoint union of sets of \mathcal{R}_J . Define $P_0(F) \equiv \nu_J(\pi_J(F))$. Then P_0 is well defined because of the consistency condition on the measures ν_J . P_0 is clearly finitely additive because the ν_J are measures and one can pick J as large as desired to include all t where there may be something other than M_t . Also, from the definition,

$$P_0(\Omega) \equiv P_0\left(\prod_{t \in I} M_t\right) = \nu_{t_1}(M_{t_1}) = 1.$$

Next I will show P_0 is a finite measure on \mathcal{E} . After this it is only a matter of using the Caratheodory extension theorem to get the existence of the desired probability measure P .

Claim: Suppose E^n is in \mathcal{E} and suppose $E^n \downarrow \emptyset$. Then $P_0(E^n) \downarrow 0$.

Proof of the claim: If not, there exists a sequence such that although $E^n \downarrow \emptyset$, $P_0(E^n) \downarrow \varepsilon > 0$. Let $E^n \in \mathcal{E}_{J_n}$. Thus it is a finite disjoint union of sets of \mathcal{R}_{J_n} . By regularity of the measures ν_{J_n} , which follows from Lemmas 9.8.4 and 9.8.5, there exists $K_{J_n} \subseteq E^n$ such that

$$\nu_{J_n}(\pi_{J_n}(K_{J_n})) + \frac{\varepsilon}{2^{n+2}} > \nu_{J_n}(\pi_{J_n}(E^n))$$

Thus $P_0(K_{J_n}) + \frac{\varepsilon}{2^{n+2}} \equiv \nu_{J_n}(\pi_{J_n}(K_{J_n})) + \frac{\varepsilon}{2^{n+2}} > \nu_{J_n}(\pi_{J_n}(E^n)) \equiv P_0(E^n)$. The interesting thing about these K_{J_n} is: they have the finite intersection property. Here is why.

$$\begin{aligned} \varepsilon &\leq P_0(\cap_{k=1}^m K_{J_k}) + P_0(E^m \setminus \cap_{k=1}^m K_{J_k}) \leq P_0(\cap_{k=1}^m K_{J_k}) + P_0(\cup_{k=1}^m E^k \setminus K_{J_k}) \\ &< P_0(\cap_{k=1}^m K_{J_k}) + \sum_{k=1}^{\infty} \frac{\varepsilon}{2^{k+2}} < P_0(\cap_{k=1}^m K_{J_k}) + \varepsilon/2, \end{aligned}$$

and so $P_0(\cap_{k=1}^m K_{J_k}) > \varepsilon/2$. In considering all the E^n , there are countably many entries in the product space which have something other than M_t in them. Say these are $\{t_1, t_2, \dots\}$.

Let p_{t_i} be a point which is in the intersection of the t_i components of the sets K_{J_n} . The compact sets in the t_i position must have the finite intersection property also because if not, the sets K_{J_n} can't have it. Thus there is such a point. As to the other positions, use the axiom of choice to pick something in each of these. Thus the intersection of these K_{J_n} contains a point which is contrary to $E^n \downarrow \emptyset$ because these sets are contained in the E^n .

With the claim, it follows P_0 is a measure on \mathcal{E} . Here is why: If $E = \bigcup_{k=1}^{\infty} E^k$ where $E, E^k \in \mathcal{E}$, then $(E \setminus \bigcup_{k=1}^n E^k) \downarrow \emptyset$ and so $P_0(\bigcup_{k=1}^n E^k) \rightarrow P_0(E)$. Hence if the E_k are disjoint, $P_0(\bigcup_{k=1}^n E_k) = \sum_{k=1}^n P_0(E_k) \rightarrow P_0(E)$. Thus for disjoint E_k having $\bigcup_k E_k = E \in \mathcal{E}$, $P_0(\bigcup_{k=1}^{\infty} E_k) = \sum_{k=1}^{\infty} P_0(E_k)$.

Now to conclude the proof, apply the Caratheodory extension theorem to obtain P a probability measure which extends P_0 to a σ algebra which contains $\sigma(\mathcal{E})$ the sigma algebra generated by \mathcal{E} with $P = P_0$ on \mathcal{E} . Thus for $E_J \in \mathcal{E}$, $P(E_J) = P_0(E_J) = v_J(P_J E_J)$.

Next, let $(\prod_{t \in I} M_t, \mathcal{F}, P)$ be the probability space and for $x \in \prod_{t \in I} M_t$ let $X_t(x) = x_t$, the t^{th} entry of x . It follows X_t is measurable (also continuous) because if U is open in M_t , then $X_t^{-1}(U)$ has a U in the t^{th} slot and M_s everywhere else for $s \neq t$. Thus inverse images of open sets are measurable. Also, letting J be a finite subset of I and for $J = (t_1, \dots, t_n)$, and F_{t_1}, \dots, F_{t_n} Borel sets in $M_{t_1} \dots M_{t_n}$ respectively, it follows F_J , where F_J has F_{t_i} in the t_i^{th} entry, is in \mathcal{E} and therefore,

$$\begin{aligned} P([X_{t_1} \in F_{t_1}] \cap [X_{t_2} \in F_{t_2}] \cap \dots \cap [X_{t_n} \in F_{t_n}]) &= \\ P([(X_{t_1}, X_{t_2}, \dots, X_{t_n}) \in F_{t_1} \times \dots \times F_{t_n}]) &= P(F_J) = P_0(F_J) \\ &= v_{t_1 \dots t_n}(F_{t_1} \times \dots \times F_{t_n}) \end{aligned}$$

Finally consider the claim about the integrals. Suppose $f(x_{t_1}, \dots, x_{t_n}) = \mathcal{X}_F$ where F is a Borel set of $\prod_{t \in J} M_t$ where $J = (t_1, \dots, t_n)$. To begin with suppose

$$F = F_{t_1} \times \dots \times F_{t_n} \quad (20.4)$$

where each F_{t_j} is in $\mathcal{B}(M_{t_j})$. Then

$$\begin{aligned} \int_{M_{t_1} \times \dots \times M_{t_n}} \mathcal{X}_F(x_{t_1}, \dots, x_{t_n}) d v_{t_1 \dots t_n} &= v_{t_1 \dots t_n}(F_{t_1} \times \dots \times F_{t_n}) \\ &= P\left(\prod_{t \in J} F_t\right) = \int_{\Omega} \mathcal{X}_{\prod_{t \in J} F_t}(x) dP = \int_{\Omega} \mathcal{X}_F(x_{t_1}, \dots, x_{t_n}) dP \end{aligned} \quad (20.5)$$

where $F_t = M_t$ if $t \notin J$. Let \mathcal{K} denote sets F of the sort in 20.4. It is clearly a π system. Now let \mathcal{G} denote those sets F in $\mathcal{B}(\prod_{t \in J} M_t)$ such that 20.5 holds. Thus $\mathcal{G} \supseteq \mathcal{K}$. It is clear that \mathcal{G} is closed with respect to countable disjoint unions and complements. Hence $\mathcal{G} \supseteq \sigma(\mathcal{K})$ but $\sigma(\mathcal{K}) = \mathcal{B}(\prod_{t \in J} M_t)$ because every open set in $\prod_{t \in J} M_t$ is the countable union of rectangles like 20.4 in which each F_{t_i} is open. Therefore, 20.5 holds for every $F \in \mathcal{B}(\prod_{t \in J} M_t)$.

Passing to simple functions and then using the monotone convergence theorem yields the final claim of the theorem. ■

As a special case, you can obtain a version of product measure for possibly infinitely many factors. Suppose in the context of the above theorem that v_t is a probability measure defined on the Borel sets of $M_t \equiv \mathbb{R}^{n_t}$ for n_t a positive integer, and let the measures, $v_{t_1 \dots t_n}$

be defined on the Borel sets of $\prod_{i=1}^n M_{t_i}$ by $v_{t_1 \dots t_n}(E) \equiv \overbrace{(v_{t_1} \times \dots \times v_{t_n})}^{\text{product measure}}(E)$. Then these

measures satisfy the necessary consistency condition and so the Kolmogorov extension theorem given above can be applied to obtain a measure P defined on a measure space $(\prod_{t \in I} M_t, \mathcal{F})$ and measurable functions $X_s : \prod_{t \in I} M_t \rightarrow M_s$ such that for F_{t_i} a Borel set in M_{t_i} ,

$$\begin{aligned} P \left((X_{t_1}, \dots, X_{t_n}) \in \prod_{i=1}^n F_{t_i} \right) &= \nu_{t_1 \dots t_n} (F_{t_1} \times \dots \times F_{t_n}) \\ &= \nu_{t_1} (F_{t_1}) \cdots \nu_{t_n} (F_{t_n}). \end{aligned} \quad (20.6)$$

In particular, $P(X_t \in F_t) = \nu_t(F_t)$. Then P in the resulting probability space given by $(\prod_{t \in I} M_t, \mathcal{F}, P)$ will be denoted as $\prod_{t \in I} \nu_t$. This proves the following theorem which describes an infinite product measure.

Theorem 20.3.4 *Let M_t for $t \in I$ be given as in Theorem 20.3.3 and let ν_t be a Borel probability measure defined on the Borel sets of M_t . Then there exists a measure P and a σ algebra $\mathcal{F} = \sigma(\mathcal{E})$ where \mathcal{E} is given in the Notation 20.3.1 such that $(\prod_t M_t, \mathcal{F}, P)$ is a probability space satisfying 20.6 whenever each F_{t_i} is a Borel set of M_{t_i} . This probability measure could be denoted as $\prod_t \nu_t$.*

20.4 Exercises

1. Suppose X and Y are metric spaces having compact closed balls. Show $(X \times Y, d_{X \times Y})$ is also a metric space which has the closures of balls compact. Here

$$d_{X \times Y}((x_1, y_1), (x_2, y_2)) \equiv \max(d(x_1, x_2), d(y_1, y_2)).$$

Let $\mathcal{A} \equiv \{E \times F : E \text{ is a Borel set in } X, F \text{ is a Borel set in } Y\}$. Show $\sigma(\mathcal{A})$, which is the smallest σ algebra containing \mathcal{A} contains the Borel sets. **Hint:** Show every open set in a metric space which has closed balls compact can be obtained as a countable union of compact sets. Next show this implies every open set can be obtained as a countable union of open sets of the form $U \times V$ where U is open in X and V is open in Y .

2. Suppose $(\Omega, \mathcal{S}, \mu)$ is a measure space which may not be complete. Could you obtain a complete measure space, $(\Omega, \mathcal{S}, \mu_1)$ by simply letting \mathcal{S} consist of all sets of the form E where there exists $F \in \mathcal{S}$ such that $(F \setminus E) \cup (E \setminus F) \subseteq N$ for some $N \in \mathcal{S}$ which has measure zero and then let $\mu(E) = \mu_1(F)$? Explain.
3. Let $(\Omega, \mathcal{S}, \mu)$ measure space and let $f : \Omega \rightarrow [0, \infty)$ be measurable. Define $A \equiv \{(x, y) : y < f(x)\}$. Show that $\int f d\mu = \int \int \mathcal{X}_A(x, y) d\mu dm$. Next show that A is product measurable in the sense that A_x is m measurable and A_y is μ measurable. Here $A_x \equiv \{y : (x, y) \in A\}$ and A_y similar. Next show that you can interchange the order of integration. **Hint:** First suppose f is a nonnegative simple function.
4. For f a nonnegative measurable function, it was shown $\int f d\mu = \int \mu([f > t]) dt$. Would it work the same if you used $\int \mu([f \geq t]) dt$? Explain.
5. Let $(\Omega, \mathcal{F}, \mu)$ be a finite measure space and suppose $\{f_n\}$ is a sequence of non-negative functions which satisfy $f_n(\omega) \leq C$ independent of n, ω . Suppose also this sequence converges to 0 in measure. That is, for all $\varepsilon > 0$, $\lim_{n \rightarrow \infty} \mu([f_n \geq \varepsilon]) = 0$. Show that then $\lim_{n \rightarrow \infty} \int_{\Omega} f_n(\omega) d\mu = 0$.

6. Explain why for each $t > 0$, $x \rightarrow e^{-tx}$ is a function in $L^1(\mathbb{R})$ and $\int_0^\infty e^{-tx} dx = \frac{1}{t}$. Thus $\int_0^R \frac{\sin(t)}{t} dt = \int_0^R \int_0^\infty \sin(t) e^{-tx} dx dt$. Now explain why you can change the order of integration in the above iterated integral. Then compute what you get. Next pass to a limit as $R \rightarrow \infty$ and show $\int_0^\infty \frac{\sin(t)}{t} dt = \frac{1}{2}\pi$.
7. Let $f(y) = g(y) = |y|^{-1/2}$ on $(-1, 0) \cup (0, 1)$ and $f(y) = g(y) = 0$ off $(-1, 0) \cup (0, 1)$. Find x where $\int_{\mathbb{R}} f(x-y)g(y)dy$ makes sense.
8. Let E_i be a Borel set in \mathbb{R} . Show that $\prod_{i=1}^n E_i$ is a Borel set in \mathbb{R}^n .
9. Let $\{a_n\}$ be an increasing sequence of numbers in $(0, 1)$ which converges to 1. Let g_n be a nonnegative function which equals zero outside (a_n, a_{n+1}) such that $\int g_n dx = 1$. Now for $(x, y) \in [0, 1] \times [0, 1]$ define $f(x, y) \equiv \sum_{k=1}^\infty g_n(y)(g_n(x) - g_{n+1}(x))$. Explain why this is actually a finite sum for each such (x, y) so there are no convergence questions in the infinite sum. Explain why f is a continuous function on $[0, 1] \times [0, 1]$. You can extend f to equal zero off $[0, 1] \times [0, 1]$ if you like. Show the iterated integrals exist but are not equal. In fact, show

$$\int_0^1 \int_0^1 f(x, y) dy dx = 1 \neq 0 = \int_0^1 \int_0^1 f(x, y) dx dy.$$

Does this example contradict the Fubini theorem, Corollary 10.14.11 on Page 307? Explain why or why not.

10. Let $f : [a, b] \rightarrow \mathbb{R}$ be Riemann integrable. Thus f is a bounded function and by Darboux's theorem, there exists a unique number between all the upper sums and lower sums of f , this number being the Riemann integral. Show that f is Lebesgue measurable and $\int_a^b f(x) dx = \int_{[a, b]} f dm$ where the second integral in the above is the Lebesgue integral taken with respect to one dimensional Lebesgue measure and the first is the ordinary Riemann integral.
11. Let $(\Omega, \mathcal{F}, \mu)$ be a σ finite measure space and let $f : \Omega \rightarrow [0, \infty)$ be measurable. Also let $\phi : [0, \infty) \rightarrow \mathbb{R}$ be increasing with $\phi(0) = 0$ and ϕ a C^1 function. Show that

$$\int_{\Omega} \phi \circ f d\mu = \int_0^\infty \phi'(t) \mu([f > t]) dt.$$

Hint: This can be done using the following steps. Let $t_i^n = i2^{-n}$. Show that

$$\mathcal{X}_{[f > t]}(\omega) = \lim_{n \rightarrow \infty} \sum_{i=0}^\infty \mathcal{X}_{[f > t_{i+1}^n]}(\omega) \mathcal{X}_{[t_i^n, t_{i+1}^n)}(t)$$

Now this is a countable sum of $\mathcal{F} \times \mathcal{B}([0, \infty))$ measurable functions and so it follows that $(t, \omega) \rightarrow \mathcal{X}_{[f > t]}(\omega)$ is $\mathcal{F} \times \mathcal{B}([0, \infty))$ measurable. Consequently, so is $\mathcal{X}_{[f > t]}(\omega) \phi(t)$. Note that it is important in the argument to have $f > t$. Now observe

$$\int_{\Omega} \phi \circ f d\mu = \int_{\Omega} \int_0^{f(\omega)} \phi'(t) dt d\mu = \int_{\Omega} \int_0^\infty \mathcal{X}_{[f > t]}(\omega) \phi'(t) dt d\mu$$

Use Fubini's theorem. For your information, this does not require the measure space to be σ finite. You can use a different argument which ties in to the first definition of the Lebesgue integral. The function $t \rightarrow \mu([f > t])$ is called the distribution function.

12. Give a different proof of the above as follows. First suppose f is a simple function, $f(\omega) = \sum_{k=1}^n a_k \chi_{E_k}(\omega)$ where the a_k are strictly increasing, $\phi(a_0) = a_0 \equiv 0$. Then explain carefully the steps to the following argument.

$$\begin{aligned}
 \int_{\Omega} \phi \circ f d\mu &= \sum_{i=1}^n \int_{\phi(a_{i-1})}^{\phi(a_i)} \mu([\phi \circ f > t]) dt = \sum_{i=1}^n \int_{\phi(a_{i-1})}^{\phi(a_i)} \sum_{k=i}^n \mu(E_k) dt \\
 &= \sum_{i=1}^n \sum_{k=i}^n \mu(E_k) \int_{a_{i-1}}^{a_i} \phi'(t) dt = \sum_{i=1}^n \int_{a_{i-1}}^{a_i} \phi'(t) \sum_{k=i}^n \mu(E_k) dt \\
 &= \sum_{i=1}^n \int_{a_{i-1}}^{a_i} \phi'(t) \mu([f > t]) dt = \int_0^{\infty} \phi'(t) \mu([f > t]) dt
 \end{aligned}$$

Note that this did not require the measure space to be σ finite and comes directly from the definition of the integral.

13. Give another argument for the above result as follows.

$$\int \phi \circ f d\mu = \int_0^{\infty} \mu([\phi \circ f > t]) dt = \int_0^{\infty} \mu([f > \phi^{-1}(t)]) dt$$

and now change the variable in the last integral, letting $\phi(s) = t$. Justify the easy manipulations.

Chapter 21

Banach Spaces

21.1 Theorems Based on Baire Category

Some examples of Banach spaces that have been discussed up to now are $\mathbb{R}^n, \mathbb{C}^n$, and $L^p(\Omega)$. Theorems about general Banach spaces are proved in this chapter. The main theorems to be presented here are the uniform boundedness theorem, the open mapping theorem, the closed graph theorem, and the Hahn Banach Theorem. The first three of these theorems come from the Baire category theorem which is about to be presented. They are topological in nature. The Hahn Banach theorem has nothing to do with topology. Banach spaces are all normed linear spaces and as such, they are all metric spaces because a normed linear space may be considered as a metric space with $d(x, y) \equiv \|x - y\|$. You can check that this satisfies all the axioms of a metric. As usual, if every Cauchy sequence converges, the metric space is called complete.

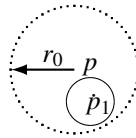
Definition 21.1.1 *A complete normed linear space is called a Banach space.*

21.1.1 Baire Category Theorem

The following remarkable result is called the Baire category theorem. To get an idea of its meaning, imagine you draw a line in the plane. The complement of this line is an open set and is dense because every point, even those on the line, are limit points of this open set. Now draw another line. The complement of the two lines is still open and dense. Keep drawing lines and looking at the complements of the union of these lines. You always have an open set which is dense. Now what if there were countably many lines? The Baire category theorem implies the complement of the union of these lines is dense. In particular it is nonempty. Thus you cannot write the plane as a countable union of lines. This is a rather rough description of this very important theorem. The precise statement and proof follow.

Theorem 21.1.2 *Let (X, d) be a complete metric space and let $\{U_n\}_{n=1}^\infty$ be a sequence of open subsets of X satisfying $\overline{U_n} = X$ (U_n is dense). Then $D \equiv \bigcap_{n=1}^\infty U_n$ is a dense subset of X .*

Proof: Let $p \in X$ and let $r_0 > 0$. I need to show $D \cap B(p, r_0) \neq \emptyset$. Since U_1 is dense, there exists $p_1 \in U_1 \cap B(p, r_0)$, an open set. Let $p_1 \in B(p_1, r_1) \subseteq \overline{B(p_1, r_1)} \subseteq U_1 \cap B(p, r_0)$ and $r_1 < 2^{-1}$. This is possible because $U_1 \cap B(p, r_0)$ is an open set and so there exists r_1 such that $B(p_1, 2r_1) \subseteq U_1 \cap B(p, r_0)$. But $B(p_1, r_1) \subseteq \overline{B(p_1, r_1)} \subseteq B(p_1, 2r_1)$ because $\overline{B(p_1, r_1)} \subseteq \{x \in X : d(x, p_1) \leq r_1\} \subseteq B(p_1, 2r_1)$. Indeed, $\{x \in X : d(x, p_1) \leq r_1\}$ is a closed set containing $B(p_1, r)$ so it contains $\overline{B(p_1, r_1)}$.



There exists $p_2 \in U_2 \cap B(p_1, r_1)$ because U_2 is dense. Let

$$p_2 \in B(p_2, r_2) \subseteq \overline{B(p_2, r_2)} \subseteq U_2 \cap B(p_1, r_1) \subseteq U_1 \cap U_2 \cap B(p, r_0).$$

and let $r_2 < 2^{-2}$. Continue in this way. Thus $r_n < 2^{-n}$,

$$\overline{B(p_n, r_n)} \subseteq U_1 \cap U_2 \cap \dots \cap U_n \cap B(p, r_0),$$

$$\overline{B(p_n, r_n)} \subseteq B(p_{n-1}, r_{n-1}).$$

The sequence, $\{p_n\}$ is a Cauchy sequence because all terms of $\{p_k\}$ for $k \geq n$ are contained in $B(p_n, r_n)$, a set whose diameter is no larger than 2^{-n} . Since X is complete, there exists p_∞ such that $\lim_{n \rightarrow \infty} p_n = p_\infty$. Since all but finitely many terms of $\{p_n\}$ are in $\overline{B(p_m, r_m)}$, it follows that $p_\infty \in \overline{B(p_m, r_m)}$ for each m . Therefore, $p_\infty \in \bigcap_{m=1}^{\infty} \overline{B(p_m, r_m)} \subseteq \bigcap_{i=1}^{\infty} U_i \cap B(p, r_0)$. ■

The following corollary is also called the Baire category theorem. It involves a countable union of closed sets rather than a countable intersection of open sets.

Corollary 21.1.3 *Let X be a complete metric space and suppose $X = \bigcup_{i=1}^{\infty} F_i$ where each F_i is a closed set. Then for some i , interior $F_i \neq \emptyset$.*

Proof: If all F_i has empty interior, then F_i^C would be a dense open set. Therefore, from Theorem 21.1.2, $\emptyset = (\bigcup_{i=1}^{\infty} F_i)^C = \bigcap_{i=1}^{\infty} F_i^C \neq \emptyset$. ■

The set D of Theorem 21.1.2 is called a G_δ set because it is the countable intersection of open sets. Thus D is a dense G_δ set.

Recall that a norm satisfies:

- a.) $\|x\| \geq 0$, $\|x\| = 0$ if and only if $x = 0$.
- b.) $\|x + y\| \leq \|x\| + \|y\|$.
- c.) $\|cx\| = |c| \|x\|$ if c is a scalar and $x \in X$.

From Theorem 3.6.2, continuity means that if $\lim_{n \rightarrow \infty} x_n = x$, then

$$\lim_{n \rightarrow \infty} f(x_n) = f(x).$$

Theorem 21.1.4 *Let X and Y be two normed linear spaces and let $L : X \rightarrow Y$ be linear ($L(ax + by) = aL(x) + bL(y)$ for a, b scalars and $x, y \in X$). The following are equivalent*

- a.) L is continuous at 0
- b.) L is continuous
- c.) There exists $K > 0$ such that $\|Lx\|_Y \leq K \|x\|_X$ for all $x \in X$ (L is bounded).

Proof: a.) \Rightarrow b.) Let $x_n \rightarrow x$. It is necessary to show that $Lx_n \rightarrow Lx$. But $(x_n - x) \rightarrow 0$ and so from continuity at 0, it follows $L(x_n - x) = Lx_n - Lx \rightarrow 0$ so $Lx_n \rightarrow Lx$. This shows a.) implies b.).

b.) \Rightarrow c.) Since L is continuous, L is continuous at 0. Hence $\|Lx\|_Y < 1$ whenever $\|x\|_X \leq \delta$ for some δ . Therefore, suppressing the subscript on the $\|$, $\|L\left(\frac{\delta x}{\|x\|}\right)\| \leq 1$. Hence $\|Lx\| \leq \frac{1}{\delta} \|x\|$.

c.) \Rightarrow a.) follows from the inequality given in c.). ■

Definition 21.1.5 *Let $L : X \rightarrow Y$ be linear and continuous where X and Y are normed linear spaces. Denote the set of all such continuous linear maps by $\mathcal{L}(X, Y)$ and define*

$$\|L\| = \sup\{\|Lx\| : \|x\| \leq 1\}. \quad (21.1)$$

This is called the operator norm.

Note that from Theorem 21.1.4 $\|L\|$ is well defined because of part c.) of that Theorem.

The next lemma follows immediately from the definition of the norm and the assumption that L is linear.

Lemma 21.1.6 *With $\|L\|$ defined in 21.1, $\mathcal{L}(X, Y)$ is a normed linear space. Also $\|Lx\| \leq \|L\| \|x\|$.*

Proof: Let $x \neq 0$ then $x/\|x\|$ has norm equal to 1 and so $\left\|L\left(\frac{x}{\|x\|}\right)\right\| \leq \|L\|$. Therefore, multiplying both sides by $\|x\|$, $\|Lx\| \leq \|L\| \|x\|$. This is obviously a linear space. It remains to verify the operator norm really is a norm. First of all, if $\|L\| = 0$, then $Lx = 0$ for all $\|x\| \leq 1$. It follows that for any $x \neq 0$, $0 = L\left(\frac{x}{\|x\|}\right)$ and so $Lx = 0$. Therefore, $L = 0$. Also, if c is a scalar,

$$\|cL\| = \sup_{\|x\| \leq 1} \|cL(x)\| = |c| \sup_{\|x\| \leq 1} \|Lx\| = |c| \|L\|.$$

It remains to verify the triangle inequality. Let $L, M \in \mathcal{L}(X, Y)$.

$$\begin{aligned} \|L + M\| &\equiv \sup_{\|x\| \leq 1} \|(L + M)(x)\| \leq \sup_{\|x\| \leq 1} (\|Lx\| + \|Mx\|) \\ &\leq \sup_{\|x\| \leq 1} \|Lx\| + \sup_{\|x\| \leq 1} \|Mx\| = \|L\| + \|M\|. \end{aligned}$$

This shows the operator norm is really a norm as hoped. ■

For example, consider the space of linear transformations defined on \mathbb{R}^n having values in \mathbb{R}^m . The fact the transformation is linear automatically imparts continuity to it. You should give a proof of this fact. Recall that every such linear transformation can be realized in terms of matrix multiplication.

Thus, in finite dimensions the algebraic condition that an operator is linear is sufficient to imply the topological condition that the operator is continuous. The situation is not so simple in infinite dimensional spaces such as $C(X; \mathbb{R}^n)$. This explains the imposition of the topological condition of continuity as a criterion for membership in $\mathcal{L}(X, Y)$ in addition to the algebraic condition of linearity. Here is an example which shows that this extra assumption cannot be eliminated.

Example 21.1.7 *Let V denote all linear combinations of functions of the form $e^{-\alpha x^2}$ for $\alpha > 0$. Thus a typical element of V is an expression of the form*

$$\sum_{k=1}^n \beta_k e^{-\alpha_k x^2}, \alpha_k > 0.$$

Let $L: V \rightarrow \mathbb{C}$ be given by $Lf \equiv \int_{\mathbb{R}} f(x) dx$. For a norm on V , $\|f\| \equiv \max\{|f(x)| : x \in \mathbb{R}\}$. Of course V is not complete, but it is a normed linear space and you could consider its completion if desired, in terms of equivalence classes of Cauchy sequences, similar to the construction of \mathbb{R} from \mathbb{Q} . Recall that $\int_{-\infty}^{\infty} e^{-x^2} dx = \int_{-\infty}^{\infty} \frac{1}{n} e^{-(x^2/n^2)} = \sqrt{\pi}$ where here $n \in \mathbb{N}$. Consider the sequence of functions $f_n(x) \equiv \frac{1}{n} e^{-(x^2/n^2)}$. Its maximum value is $1/n$ and so $\|f_n\| \rightarrow 0$ but Lf_n fails to converge to 0. Thus L is not continuous although it is linear.

Theorem 21.1.8 *If Y is a Banach space, then $\mathcal{L}(X, Y)$ is also a Banach space.*

Proof: Let $\{L_n\}$ be a Cauchy sequence in $\mathcal{L}(X, Y)$ and let $x \in X$.

$$\|L_n x - L_m x\| \leq \|x\| \|L_n - L_m\|.$$

Thus $\{L_n x\}$ is a Cauchy sequence. Let $Lx = \lim_{n \rightarrow \infty} L_n x$. Then, clearly, L is linear because if x_1, x_2 are in X , and a, b are scalars, then

$$L(ax_1 + bx_2) = \lim_{n \rightarrow \infty} L_n(ax_1 + bx_2) = \lim_{n \rightarrow \infty} (aL_n x_1 + bL_n x_2) = aLx_1 + bLx_2.$$

Also L is continuous. To see this, note that $\{\|L_n\|\}$ is a Cauchy sequence of real numbers because $\|\|L_n\| - \|L_m\|\| \leq \|L_n - L_m\|$. Hence there exists $K > \sup\{\|L_n\| : n \in \mathbb{N}\}$. Thus, if $x \in X$, $\|Lx\| = \lim_{n \rightarrow \infty} \|L_n x\| \leq K \|x\|$. ■

21.1.2 Uniform Boundedness Theorem

The next big result is sometimes called the Uniform Boundedness theorem, or the Banach-Steinhaus theorem. This is a very surprising theorem which implies that for a collection of bounded linear operators, if they are bounded pointwise, then they are also bounded uniformly. As an example of a situation in which pointwise bounded does not imply uniformly bounded, consider the functions $f_\alpha(x) \equiv \mathcal{X}_{(\alpha, 1)}(x)x^{-1}$ for $\alpha \in (0, 1)$. Clearly each function is bounded and the collection of functions is bounded at each point of $(0, 1)$, but there is no bound for all these functions taken together.

Theorem 21.1.9 *Let X be a Banach space and let Y be a normed linear space. Let $\{L_\alpha\}_{\alpha \in \Lambda}$ be a collection of elements of $\mathcal{L}(X, Y)$. Then one of the following happens.*

- a.) $\sup\{\|L_\alpha\| : \alpha \in \Lambda\} < \infty$
- b.) *There exists a dense G_δ set D , such that for all $x \in D$,*

$$\sup\{\|L_\alpha x\| : \alpha \in \Lambda\} = \infty.$$

Proof: For each $n \in \mathbb{N}$, define $U_n = \{x \in X : \sup\{\|L_\alpha x\| : \alpha \in \Lambda\} > n\}$. Then U_n is an open set because if $x \in U_n$, then there exists $\alpha \in \Lambda$ such that $\|L_\alpha x\| > n$. But then, since L_α is continuous, this situation persists for all y sufficiently close to x , say for all $y \in B(x, \delta)$. Then $B(x, \delta) \subseteq U_n$ which shows U_n is open.

Case b.) is obtained from Theorem 21.1.2 if each U_n is dense.

The other case is that for some n , U_n is not dense. If this occurs, there exists x_0 and $r > 0$ such that for all $B(x_0, r) \subseteq U_n^c$ so for all $x \in B(x_0, r)$, $\|L_\alpha x\| \leq n$ for all α . Now if $y \in B(0, r)$, $x_0 + y \in B(x_0, r)$. Consequently, for all such y , $\|L_\alpha(x_0 + y)\| \leq n$. This implies that for all $\alpha \in \Lambda$ and $\|y\| < r$, $\|L_\alpha y\| \leq n + \|L_\alpha(x_0)\| \leq 2n$. Therefore, if $\|y\| \leq 1$, $\|\frac{r}{2}y\| < r$ and so for all α , $\|L_\alpha(\frac{r}{2}y)\| \leq 2n$. Now multiplying by $r/2$ it follows that whenever $\|y\| \leq 1$, $\|L_\alpha(y)\| \leq 4n/r$. Hence case a.) holds. ■

21.1.3 Open Mapping Theorem

Another remarkable theorem which depends on the Baire category theorem is the open mapping theorem. Unlike Theorem 21.1.9 it requires both X and Y to be Banach spaces.

Theorem 21.1.10 *Let X and Y be Banach spaces, let $L \in \mathcal{L}(X, Y)$, and suppose L is onto. Then L maps open sets onto open sets.*

To aid in the proof, here is a lemma.

Lemma 21.1.11 *Let a and b be positive constants and suppose $B(0, a) \subseteq \overline{L(B(0, b))}$. Then $\overline{L(B(0, b))} \subseteq L(B(0, 2b))$.*

Proof: Let $z \in \overline{L(B(0, b))}$. Then let $x_1 \in B(0, b)$ with $\|z - Lx_1\| < \frac{a}{2^1}$. Then it follows that on multiplying by 2, $\|2z - 2Lx_1\| < a$ and so there is $x_2 \in B(0, b)$ with

$$\|(2z - 2Lx_1) - Lx_2\| < \frac{a}{2}$$

which implies $\left\|z - \left(Lx_1 + \frac{Lx_2}{2}\right)\right\| < \frac{a}{2^2}$. If $x_i \in B(0, b)$ and $\left\|z - \sum_{i=0}^n \frac{Lx_i}{2^i}\right\| < \frac{a}{2^n}$, then

$$\left\|2^n \left(z - \sum_{i=0}^n \frac{Lx_i}{2^i}\right)\right\| < a$$

Continuing this way, there is $x_{n+1} \in B(0, b)$ such that

$$\left\|2^n \left(z - \sum_{i=0}^n \frac{Lx_i}{2^i}\right) - L(x_{n+1})\right\| < \frac{a}{2}$$

and so $\left\|z - \sum_{i=0}^{n+1} \frac{Lx_i}{2^i}\right\| < \frac{a}{2^{n+1}}$. Let $x \equiv \sum_{i=0}^{\infty} \frac{x_i}{2^i}$. Then from the triangle inequality, $\|x\| < \sum_{i=0}^{\infty} \frac{b}{2^i} = 2b$ and by continuity of L ,

$$\|z - Lx\| = \lim_{n \rightarrow \infty} \left\|z - L\left(\sum_{i=0}^n \frac{x_i}{2^i}\right)\right\| \leq \lim_{n \rightarrow \infty} \frac{a}{2^n} = 0 \blacksquare$$

Proof of Theorem 21.1.10: $Y = \cup_{n=1}^{\infty} \overline{L(B(0, n))}$. By Corollary 21.1.3, $\overline{L(B(0, n_0))}$ has nonempty interior for some n_0 . Thus $B(y, r) \subseteq \overline{L(B(0, n_0))}$ for some y and some $r > 0$. Since L is linear $B(-y, r) \subseteq \overline{L(B(0, n_0))}$ also. Here is why. If $z \in B(-y, r)$, then $-z \in B(y, r)$ and so there exists $x_n \in B(0, n_0)$ such that $Lx_n \rightarrow -z$. Therefore, $L(-x_n) \rightarrow z$ and $-x_n \in B(0, n_0)$ also. Therefore $z \in \overline{L(B(0, n_0))}$. Then it follows that

$$\begin{aligned} B(0, r) &\subseteq B(y, r) + B(-y, r) \equiv \{y_1 + y_2 : y_1 \in B(y, r) \text{ and } y_2 \in B(-y, r)\} \\ &\subseteq \overline{L(B(0, 2n_0))} \end{aligned}$$

The reason for the last inclusion is that from the above, if $y_1 \in B(y, r)$ and $y_2 \in B(-y, r)$, there exists $x_n, z_n \in B(0, n_0)$ such that $Lx_n \rightarrow y_1$, $Lz_n \rightarrow y_2$. Therefore, $\|x_n + z_n\| \leq 2n_0$ and so $(y_1 + y_2) \in \overline{L(B(0, 2n_0))}$.

By Lemma 21.1.11, $\overline{L(B(0, 2n_0))} \subseteq L(B(0, 4n_0))$ so $B(0, r) \subseteq L(B(0, 4n_0))$. Letting $a = r(4n_0)^{-1}$, it follows, since L is linear, that $B(0, a) \subseteq L(B(0, 1))$. It follows since L is linear,

$$L(B(0, r)) \supseteq B(0, ar). \quad (21.2)$$

Now let U be open in X and let $x + B(0, r) = B(x, r) \subseteq U$. Using 21.2,

$$\begin{aligned} L(U) &\supseteq L(x + B(0, r)) \\ &= Lx + L(B(0, r)) \supseteq Lx + B(0, ar) = B(Lx, ar). \end{aligned}$$

Hence $Lx \in B(Lx, ar) \subseteq L(U)$ which shows that every point, $Lx \in LU$, is an interior point of LU and so LU is open. ■

This theorem is surprising because it implies that if $|\cdot|$ and $\|\cdot\|$ are two norms with respect to which a vector space X is a Banach space such that $|\cdot| \leq K \|\cdot\|$, then there exists a constant k , such that $\|\cdot\| \leq k |\cdot|$. This can be useful because sometimes it is not clear how to compute k when all that is needed is its existence. To see the open mapping theorem implies this, consider the identity map $\text{id}x = x$. Then $\text{id} : (X, \|\cdot\|) \rightarrow (X, |\cdot|)$ is continuous and onto. Hence id is an open map which implies id^{-1} is continuous. Theorem 21.1.4 gives the existence of the constant k .

21.1.4 Closed Graph Theorem

Definition 21.1.12 Let $f : D \rightarrow E$. The set of all ordered pairs of the form

$$\{(x, f(x)) : x \in D\}$$

is called the graph of f .

Definition 21.1.13 If X and Y are normed linear spaces, make $X \times Y$ into a normed linear space by using the norm $\|(x, y)\| = \max(\|x\|, \|y\|)$ along with component-wise addition and scalar multiplication. Thus $a(x, y) + b(z, w) \equiv (ax + bz, ay + bw)$.

There are other ways to give a norm for $X \times Y$. For example, you could define $\|(x, y)\| = \|x\| + \|y\|$

Lemma 21.1.14 The norm defined in Definition 21.1.13 on $X \times Y$ along with the definition of addition and scalar multiplication given there make $X \times Y$ into a normed linear space.

Proof: The only axiom for a norm which is not obvious is the triangle inequality. Therefore, consider

$$\begin{aligned} \|(x_1, y_1) + (x_2, y_2)\| &= \|(x_1 + x_2, y_1 + y_2)\| \\ &= \max(\|x_1 + x_2\|, \|y_1 + y_2\|) \\ &\leq \max(\|x_1\| + \|x_2\|, \|y_1\| + \|y_2\|) \end{aligned}$$

Both $\|x_1\| + \|x_2\|$ and $\|y_1\| + \|y_2\|$ are no larger than

$$\max(\|x_1\|, \|y_1\|) + \max(\|x_2\|, \|y_2\|)$$

and so the above is

$$\leq \max(\|x_1\|, \|y_1\|) + \max(\|x_2\|, \|y_2\|) = \|(x_1, y_1)\| + \|(x_2, y_2)\|.$$

It is obvious $X \times Y$ is a vector space from the above definition. ■

Lemma 21.1.15 If X and Y are Banach spaces, then $X \times Y$ with the norm and vector space operations defined in Definition 21.1.13 is also a Banach space.

Proof: The only thing left to check is that the space is complete. But this follows from the simple observation that $\{(x_n, y_n)\}$ is a Cauchy sequence in $X \times Y$ if and only if $\{x_n\}$ and $\{y_n\}$ are Cauchy sequences in X and Y respectively. Thus if $\{(x_n, y_n)\}$ is a Cauchy sequence in $X \times Y$, it follows there exist x and y such that $x_n \rightarrow x$ and $y_n \rightarrow y$. But then from the definition of the norm, $(x_n, y_n) \rightarrow (x, y)$. ■

Lemma 21.1.16 *Every closed subspace of a Banach space is a Banach space.*

Proof: If $F \subseteq X$ where X is a Banach space and $\{x_n\}$ is a Cauchy sequence in F , then since X is complete, there exists a unique $x \in X$ such that $x_n \rightarrow x$. However this means $x \in \overline{F} = F$ since F is closed. ■

Definition 21.1.17 *Let X and Y be Banach spaces and let $D \subseteq X$ be a subspace. A linear map $L : D \rightarrow Y$ is said to be closed if its graph is a closed subspace of $X \times Y$. Equivalently, L is closed if $x_n \rightarrow x$ and $Lx_n \rightarrow y$ implies $x \in D$ and $y = Lx$.*

Note the distinction between closed and continuous. If the operator is closed the assertion that $y = Lx$ only follows if it is known that the sequence $\{Lx_n\}$ converges. In the case of a continuous operator, the convergence of $\{Lx_n\}$ follows from the assumption that $x_n \rightarrow x$. It is not always the case that a mapping which is closed is necessarily continuous. Consider the function $f(x) = \tan(x)$ if x is not an odd multiple of $\frac{\pi}{2}$ and $f(x) \equiv 0$ at every odd multiple of $\frac{\pi}{2}$. Then the graph is closed and the function is defined on \mathbb{R} but it clearly fails to be continuous. Of course this function is not linear. You could also consider the map,

$$\frac{d}{dx} : \{y \in C^1([0, 1]) : y(0) = 0\} \equiv D \rightarrow C([0, 1]).$$

where the norm is the uniform norm on $C([0, 1])$, $\|y\|_\infty$. If $y \in D$, then $y(x) = \int_0^x y'(t) dt$. Therefore, if $\frac{dy_n}{dx} \rightarrow f \in C([0, 1])$ and if $y_n \rightarrow y$ in $C([0, 1])$ it follows that

$$\begin{array}{ccc} y_n(x) & = & \int_0^x \frac{dy_n(t)}{dx} dt \\ \downarrow & & \downarrow \\ y(x) & = & \int_0^x f(t) dt \end{array}$$

and so by the fundamental theorem of calculus $f(x) = y'(x)$ and so the mapping is closed. It is obviously not continuous because it takes $y(x)$ and $y(x) + \frac{1}{n} \sin(nx)$ to two functions which are far from each other even though these two functions are very close in $C([0, 1])$. Furthermore, it is not defined on the whole space, $C([0, 1])$.

The next theorem, the closed graph theorem, gives conditions under which closed implies continuous.

Theorem 21.1.18 *Let X and Y be Banach spaces and suppose $L : X \rightarrow Y$ is closed and linear. Then L is continuous.*

Proof: Let G be the graph of L . $G = \{(x, Lx) : x \in X\}$. By Lemma 21.1.16 it follows that G is a Banach space. Define $P : G \rightarrow X$ by $P(x, Lx) = x$. P maps the Banach space G onto the Banach space X and is continuous and linear. By the open mapping theorem, P maps open sets onto open sets. Since P is also one to one, this says that P^{-1} is continuous. Thus $\|P^{-1}x\| \leq K\|x\|$. Hence

$$\|Lx\| \leq \max(\|x\|, \|Lx\|) \leq K\|x\|$$

By Theorem 21.1.4 on Page 534, this shows L is continuous. ■

The following corollary is quite useful. It shows how to obtain a new norm on the domain of a closed operator such that the domain with this new norm becomes a Banach space.

Corollary 21.1.19 *Let $L : D \subseteq X \rightarrow Y$ where X, Y are a Banach spaces, and L is a closed operator. Then define a new norm on D by*

$$\|x\|_D \equiv \|x\|_X + \|Lx\|_Y.$$

Then D with this new norm is a Banach space.

Proof: If $\{x_n\}$ is a Cauchy sequence in D with this new norm, it follows both $\{x_n\}$ and $\{Lx_n\}$ are Cauchy sequences and therefore, they converge. Since L is closed, $x_n \rightarrow x$ and $Lx_n \rightarrow Lx$ for some $x \in D$. Thus $\|x_n - x\|_D \rightarrow 0$. ■

21.2 Hahn Banach Theorem

The closed graph, open mapping, and uniform boundedness theorems are the three major topological theorems in functional analysis. The other major theorem is the Hahn-Banach theorem which has nothing to do with topology. Before presenting this theorem, here are some preliminaries about partially ordered sets.

21.2.1 Partially Ordered Sets

Recall Theorem 2.8.2 which is stated next for convenience.

Theorem 21.2.1 (Hausdorff Maximal Principle) *Let \mathcal{F} be a nonempty partially ordered set. Then there exists a maximal chain.*

21.2.2 Gauge Functions and Hahn Banach Theorem

Definition 21.2.2 *Let X be a real vector space $\rho : X \rightarrow \mathbb{R}$ is called a gauge function if*

$$\begin{aligned} \rho(x+y) &\leq \rho(x) + \rho(y), \\ \rho(ax) &= a\rho(x) \text{ if } a \geq 0. \end{aligned} \tag{21.3}$$

Suppose M is a subspace of X and $z \notin M$. Suppose also that f is a linear real-valued function having the property that $f(x) \leq \rho(x)$ for all $x \in M$. Consider the problem of extending f to $M \oplus \mathbb{R}z$ such that if F is the extended function, $F(y) \leq \rho(y)$ for all $y \in M \oplus \mathbb{R}z$ and F is linear. Since F is to be linear, it suffices to determine how to define $F(z)$. Letting $a > 0$, it is required to define $F(z)$ such that the following hold for all $x, y \in M$.

$$\begin{aligned} \overbrace{F(x)}^{f(x)} + aF(z) &= F(x+az) \leq \rho(x+az), \\ \overbrace{F(y)}^{f(y)} - aF(z) &= F(y-az) \leq \rho(y-az). \end{aligned} \tag{21.4}$$

Note that something in $M \oplus \mathbb{R}z$ is of the form $x + az$ or $x - az$ for $a \geq 0$ and this includes all possibilities. Now multiplying by a^{-1} for $a > 0$, the above holds if and only if for all $x, y \in M$,

$$F(y) - F(z) \leq \rho(y - z), F(x) + F(z) \leq \rho(x + z)$$

Thus we need to choose $F(z)$ such that for all $x, y \in M$,

$$f(y) - \rho(y - z) \leq F(z) \leq \rho(x + z) - f(x). \quad (21.5)$$

Is there any such number between $f(y) - \rho(y - z)$ and $\rho(x + z) - f(x)$ for every pair $x, y \in M$? This is where $f(x) \leq \rho(x)$ on M and that f is linear is used. For $x, y \in M$,

$$\begin{aligned} & \rho(x + z) - f(x) - [f(y) - \rho(y - z)] \\ &= \rho(x + z) + \rho(y - z) - (f(x) + f(y)) \\ &\geq \rho(x + y) - f(x + y) \geq 0. \end{aligned}$$

Then if $a = \sup \{f(y) - \rho(y - z) : y \in M\}$, $b = \inf \{\rho(x + z) - f(x) : x \in M\}$, it follows that $[a, b] \neq \emptyset$. Choose $F(z)$ in $[a, b]$ so it will satisfy 21.5. This has proved the following lemma.

Lemma 21.2.3 *Let M be a subspace of X , a real linear space, and let ρ be a gauge function on X . Suppose $f : M \rightarrow \mathbb{R}$ is linear, $z \notin M$, and $f(x) \leq \rho(x)$ for all $x \in M$. Then f can be extended to $M \oplus \mathbb{R}z$ such that, if F is the extended function, F is linear and $F(x) \leq \rho(x)$ for all $x \in M \oplus \mathbb{R}z$.*

With this lemma, the Hahn Banach theorem can be proved.

Theorem 21.2.4 (Hahn Banach theorem) *Let X be a real vector space, let M be a subspace of X , let $f : M \rightarrow \mathbb{R}$ be linear, let ρ be a gauge function on X , and suppose $f(x) \leq \rho(x)$ for all $x \in M$. Then there exists a linear function, $F : X \rightarrow \mathbb{R}$, such that*

- a.) $F(x) = f(x)$ for all $x \in M$
- b.) $F(x) \leq \rho(x)$ for all $x \in X$.

Proof: Let $\mathcal{F} = \{(V, g) : V \supseteq M, V \text{ is a subspace of } X, g : V \rightarrow \mathbb{R} \text{ is linear, } g(x) = f(x) \text{ for all } x \in M, \text{ and } g(x) \leq \rho(x) \text{ for } x \in V\}$. Then $(M, f) \in \mathcal{F}$ so $\mathcal{F} \neq \emptyset$. Define a partial order by the following rule. $(V, g) \leq (W, h)$ means

$$V \subseteq W \text{ and } h(x) = g(x) \text{ if } x \in V.$$

By Theorem 21.2.1, there exists a maximal chain, $\mathcal{C} \subseteq \mathcal{F}$. Let $Y = \cup \{V : (V, g) \in \mathcal{C}\}$ and let $h : Y \rightarrow \mathbb{R}$ be defined by $h(x) = g(x)$ where $x \in V$ and $(V, g) \in \mathcal{C}$. This is well defined because if $x \in V_1$ and V_2 where (V_1, g_1) and (V_2, g_2) are both in the chain, then since \mathcal{C} is a chain, the two elements are related. Therefore, $g_1(x) = g_2(x)$. Also h is linear because if $ax + by \in Y$, then $x \in V_1$ and $y \in V_2$ where (V_1, g_1) and (V_2, g_2) are elements of \mathcal{C} . Therefore, letting V denote the larger of the two V_i , and g be the function that goes with V , it follows $ax + by \in V$ where $(V, g) \in \mathcal{C}$. Therefore,

$$h(ax + by) = g(ax + by) = ag(x) + bg(y) = ah(x) + bh(y).$$

Also, $h(x) = g(x) \leq \rho(x)$ for any $x \in Y$ because for such x , $x \in V$ where $(V, g) \in \mathcal{C}$.

Is $Y = X$? If not, there exists $z \in X \setminus Y$ and there exists an extension of h to $Y \oplus \mathbb{R}z$ using Lemma 21.2.3. Letting \bar{h} denote this extended function, contradicts the maximality of \mathcal{C} . Indeed, $\mathcal{C} \cup \{(Y \oplus \mathbb{R}z, \bar{h})\}$ would be a longer chain. ■

This is the original version of the theorem. There is also a version of this theorem for complex vector spaces which is based on a trick.

21.2.3 The Complex Version of the Hahn Banach Theorem

First is a lemma which is quite interesting for its own sake.

Lemma 21.2.5 *Let $h : V \rightarrow \mathbb{R}$ where V is a complex normed linear space. Then h is linear with respect to real scalars if and only if $F(x) \equiv h(x) - ih(ix)$ is linear with respect to complex scalars.*

Proof: \Leftarrow By assumption, F is linear with respect to real scalars. Let $c \in \mathbb{R}$. Then

$$cF(x) = ch(x) - cih(ix) = F(cx) = h(cx) - ih(cix).$$

Equating real parts, it follows $ch(x) = h(cx)$. Also

$$\begin{aligned} F(x+y) &= h(x+y) - ih(i(x+y)) = F(x) + F(y) \\ &= h(x) - ih(ix) + h(y) - ih(iy) \end{aligned}$$

Equating real parts, $h(x+y) = h(x) + h(y)$.

\Rightarrow I need to show that F is linear with respect to complex scalars.

$$\begin{aligned} F(ix) &\equiv h(ix) - ih(-x) = h(ix) + ih(x) \\ &= i(h(x) - ih(ix)) = iF(x) \end{aligned}$$

It is fairly obvious that $F(x+y) = F(x) + F(y)$. Also, if c is real, it is clear that $F(cx) = cF(x)$. Therefore,

$$\begin{aligned} F((a+ib)x) &= F(ax) + F(ibx) \\ &= aF(x) + ibF(x) = (a+ib)F(x) \quad \blacksquare \end{aligned}$$

Corollary 21.2.6 (*Hahn Banach*) *Let M be a subspace of a complex normed linear space X , and suppose $f : M \rightarrow \mathbb{C}$ is linear and satisfies $|f(x)| \leq K\|x\|$ for all $x \in M$. Then there exists a linear function F , defined on all of X such that $F(x) = f(x)$ for all $x \in M$ and $|F(x)| \leq K\|x\|$ for all $x \in X$.*

Proof: First note $f(x) = \operatorname{Re} f(x) + i\operatorname{Im} f(x)$ and so

$$\operatorname{Re} f(ix) + i\operatorname{Im} f(ix) = f(ix) = if(x) = i\operatorname{Re} f(x) - \operatorname{Im} f(x).$$

Therefore, $\operatorname{Im} f(x) = -\operatorname{Re} f(ix)$, and

$$f(x) = \operatorname{Re} f(x) - i\operatorname{Re} f(ix).$$

This is important because it shows it is only necessary to consider $\operatorname{Re} f$ in understanding f . From Lemma 21.2.5 $\operatorname{Re} f$ is linear with respect to real scalars.

Consider X as a real vector space and let $\rho(x) \equiv K\|x\|$. Then for all $x \in M$,

$$|\operatorname{Re} f(x)| \leq |f(x)| \leq K\|x\| \equiv \rho(x).$$

From Theorem 21.2.4, $\operatorname{Re} f$ may be extended to a function h which satisfies

$$\begin{aligned} h(ax+by) &= ah(x) + bh(y) \text{ if } a, b \in \mathbb{R} \\ h(x) &\leq K\|x\| \text{ for all } x \in X. \end{aligned}$$

Actually, $|h(x)| \leq K \|x\|$. The reason for this is that $h(-x) = -h(x) \leq K \|-x\| = K \|x\|$ and therefore, $h(x) \geq -K \|x\|$ so $-h(x) \leq K \|x\|$. Thus $|h(x)| \leq K \|x\|$. Let $F(x) \equiv h(x) - ih(ix)$. By Lemma 21.2.5, F is complex linear.

Now $wF(x) = |F(x)|$ for some $|w| = 1$. Therefore

$$\begin{aligned} |F(x)| &= wF(x) = F(wx) \equiv h(wx) - \overbrace{ih(iwx)}^{\text{must equal zero}} = h(wx) \\ &= |h(wx)| \leq K \|wx\| = K \|x\|. \blacksquare \end{aligned}$$

21.2.4 The Dual Space and Adjoint Operators

Definition 21.2.7 Let X be a Banach space. Denote by X' the space of continuous linear functions which map X to the field of scalars. Thus $X' = \mathcal{L}(X, \mathbb{F})$. By Theorem 21.1.8 on Page 535, X' is a Banach space. Remember with the norm defined on $\mathcal{L}(X, \mathbb{F})$,

$$\|f\| = \sup\{|f(x)| : \|x\| \leq 1\}$$

X' is called the dual space.

Definition 21.2.8 Let X and Y be Banach spaces and suppose $L \in \mathcal{L}(X, Y)$. Then define the adjoint map in $\mathcal{L}(Y', X')$, denoted by L^* , by

$$L^*y^*(x) \equiv y^*(Lx)$$

for all $y^* \in Y'$.

The following diagram is a good one to help remember this definition.

$$\begin{array}{ccc} & L^* & \\ X' & \leftarrow & Y' \\ & \rightarrow & \\ X & \xrightarrow{L} & Y \end{array}$$

This is a generalization of the adjoint of a linear transformation on an inner product space from Linear Algebra. Recall

$$(Ax, y) = (x, A^*y)$$

What is being done here is to generalize this algebraic concept to arbitrary Banach spaces. There are some issues which need to be discussed relative to the above definition. First of all, it must be shown that $L^*y^* \in X'$. Also, it will be useful to have the following lemma which is a useful application of the Hahn Banach theorem.

Lemma 21.2.9 Let X be a normed linear space and let $x \in X \setminus V$ where V is a closed subspace of X . Then there exists $x^* \in X'$ such that $x^*(x) = \|x\| \neq 0$, $x^*(V) = \{0\}$, and $\|x^*\| = \frac{1}{\text{dist}(x, V)} \|x\|$. In the case that $V = \{0\}$, $\|x^*\| = 1$.

Proof: Let $f : \mathbb{F}x + V \rightarrow \mathbb{F}$ be defined by $f(\alpha x + v) = \alpha \|x\|$. First it is necessary to show f is well defined and continuous. If $\alpha_1 x + v_1 = \alpha_2 x + v_2$ then if $\alpha_1 \neq \alpha_2$, then $x \in V$ which is assumed not to happen so f is well defined. It remains to show f is continuous. Suppose

then that $\alpha_n x + v_n \rightarrow 0$. It is necessary to show $\alpha_n \rightarrow 0$. If this does not happen, then there exists a subsequence, still denoted by α_n such that $|\alpha_n| \geq \delta > 0$. Then $x + (1/\alpha_n)v_n \rightarrow 0$. Thus $\|x - (1/\alpha_n)v_n\| \rightarrow 0$ so x is a limit of points of V contradicting the assumption that $x \notin V$ and V is a closed subspace. Hence f is continuous on $\mathbb{F}x + V$. Thus

$$\|f\| = \sup_{\|\alpha x + v\| \leq 1} |\alpha| \|x\|.$$

What is $\sup |\alpha|$, given that $\|\alpha x + v\| = |\alpha| \|x + \frac{v}{\alpha}\| \leq 1$? Since $\frac{v}{\alpha}$ is a generic element of V this reduces to $\sup |\alpha|$ such that $|\alpha| \|x + v\| \leq 1$. Thus the largest $|\alpha|$ can be is $\frac{1}{\text{dist}(x, V)}$ and so

$$\|f\| = \frac{1}{\text{dist}(x, V)} \|x\|$$

Now for $z \in \mathbb{F}x + V$, $|f(z)| \leq \|f\| \|z\| = \frac{1}{\text{dist}(x, V)} \|x\| \|z\|$. By the complex Hahn Banach theorem, there exists $x^* \in X'$ such that $x^* = f$ on $\mathbb{F}x + V$ and for all $z \in X$,

$$|x^*(z)| \leq \|f\| \|z\| = \frac{1}{\text{dist}(x, V)} \|x\| \|z\|$$

Thus $\|x^*\| \leq \|f\|$. However, equality must occur because

$$\begin{aligned} \frac{1}{\text{dist}(x, V)} \|x\| &= \|f\| \equiv \sup_{\|z\| \leq 1, z \in \mathbb{F}x + V} |f(z)| \\ &= \sup_{\|z\| \leq 1, z \in \mathbb{F}x + V} |x^*(z)| \leq \sup_{\|z\| \leq 1} |x^*(z)| \equiv \|x^*\|. \end{aligned}$$

In case $V = \{0\}$, $\text{dist}(x, V) = \|x\|$ and so $\|x^*\| = 1$. ■

Theorem 21.2.10 *Let $L \in \mathcal{L}(X, Y)$ where X and Y are Banach spaces. Then*

- a.) $L^* \in \mathcal{L}(Y', X')$ as claimed and $\|L^*\| = \|L\|$.*
- b.) If L maps one to one onto a closed subspace of Y , then L^* is onto.*
- c.) If L maps onto a dense subset of Y , then L^* is one to one.*

Proof: It is routine to verify L^*y^* and L^* are both linear. This follows immediately from the definition. As usual, the interesting thing concerns continuity.

$$\|L^*y^*\| = \sup_{\|x\| \leq 1} |L^*y^*(x)| = \sup_{\|x\| \leq 1} |y^*(Lx)| \leq \|y^*\| \|L\|.$$

Thus L^* is continuous as claimed and $\|L^*\| \leq \|L\|$.

By Lemma 21.2.9, there exists $y_x^* \in Y'$ such that $\|y_x^*\| = 1$ and $y_x^*(Lx) = \|Lx\|$. Therefore,

$$\begin{aligned} \|L^*\| &= \sup_{\|y^*\| \leq 1} \|L^*y^*\| = \sup_{\|y^*\| \leq 1} \sup_{\|x\| \leq 1} |L^*y^*(x)| \\ &= \sup_{\|y^*\| \leq 1} \sup_{\|x\| \leq 1} |y^*(Lx)| = \sup_{\|x\| \leq 1} \sup_{\|y^*\| \leq 1} |y^*(Lx)| \\ &\geq \sup_{\|x\| \leq 1} |y_x^*(Lx)| = \sup_{\|x\| \leq 1} \|Lx\| = \|L\| \end{aligned}$$

showing that $\|L^*\| \geq \|L\|$ and this shows part a.).

Next consider b.). Let $x^* \in X'$. Is there $y^* \in Y'$ such that $L^*(y^*) = x^*$? This will be so if and only if for all $x \in X$, $y^*(Lx) = x^*(x)$. Let $f(Lx) \equiv x^*(x)$. This defines f on $L(X)$ because of the assumption that L is one to one. Is f continuous on $L(X)$? Suppose $Lx_n \rightarrow Lx$ in $L(X)$. Does it follow that $x_n \rightarrow x$? Yes, because $L(X)$ is a closed subspace of a Banach space Y and is therefore also a Banach space. Since L is one to one, it follows from the open mapping theorem that L^{-1} is continuous on $L(X)$ and so indeed, $x_n \rightarrow x$ and so $x^*(x_n) \rightarrow x^*(x)$ showing that f is continuous on $L(X)$. Now by the Hahn Banach theorem, there is an extension y^* of f to all of Y which has the same norm. Thus L^* is onto as claimed.

Consider the last assertion. Suppose $L^*y^* = 0$. Is $y^* = 0$? Letting $Lx \in D$ where D is the dense subset of Y , $y^*(Lx) = L^*y^*(x) = 0$ and so y^* sends all in a dense subset of Y to 0. Hence, by continuity of y^* , it equals 0. Thus L^* is one to one. ■

Corollary 21.2.11 *Suppose X and Y are Banach spaces, $L \in \mathcal{L}(X, Y)$, and L is one to one and onto. Then L^* is also one to one and onto.*

There exists a natural mapping, called the James map from a normed linear space X , to the dual of the dual space which is described in the following definition.

Definition 21.2.12 *Define $J : X \rightarrow X''$ by $J(x)(x^*) = x^*(x)$.*

Theorem 21.2.13 *The map J has the following properties.*

- a.) J is one to one and linear.
 - b.) $\|Jx\| = \|x\|$ and $\|J\| = 1$.
 - c.) $J(X)$ is a closed subspace of X'' if X is complete.
- Also if $x^* \in X'$,

$$\|x^*\| = \sup \{ |x^{**}(x^*)| : \|x^{**}\| \leq 1, x^{**} \in X'' \}.$$

Proof:

$$\begin{aligned} J(ax + by)(x^*) &\equiv x^*(ax + by) = ax^*(x) + bx^*(y) \\ &= (aJ(x) + bJ(y))(x^*). \end{aligned}$$

Since this holds for all $x^* \in X'$, it follows that $J(ax + by) = aJ(x) + bJ(y)$ and so J is linear. If $Jx = 0$, then by Lemma 21.2.9 there exists x^* such that $x^*(x) = \|x\|$ and $\|x^*\| = 1$. Then $0 = J(x)(x^*) = x^*(x) = \|x\|$. This shows a.).

To show b.), let $x \in X$ and use Lemma 21.2.9 to obtain $x^* \in X'$ such that $x^*(x) = \|x\|$ with $\|x^*\| = 1$. Then

$$\begin{aligned} \|x\| &\geq \sup \{ |y^*(x)| : \|y^*\| \leq 1 \} = \sup \{ |J(x)(y^*)| : \|y^*\| \leq 1 \} = \|Jx\| \\ &\geq |J(x)(x^*)| = |x^*(x)| = \|x\| \end{aligned}$$

Therefore, $\|Jx\| = \|x\|$ as claimed. Therefore,

$$\|J\| = \sup \{ \|Jx\| : \|x\| \leq 1 \} = \sup \{ \|x\| : \|x\| \leq 1 \} = 1.$$

This shows b.).

To verify c.), use b.). If $Jx_n \rightarrow y^{**} \in X''$ then by b.), x_n is a Cauchy sequence converging to some $x \in X$ because $\|x_n - x_m\| = \|Jx_n - Jx_m\|$ and $\{Jx_n\}$ is a Cauchy sequence. Then $Jx = \lim_{n \rightarrow \infty} Jx_n = y^{**}$.

Finally, to show the assertion about the norm of x^* , use what was just shown applied to the James map from X' to X''' still referred to as J .

$$\begin{aligned}\|x^*\| &= \sup \{|x^*(x)| : \|x\| \leq 1\} = \sup \{|J(x)(x^*)| : \|Jx\| \leq 1\} \\ &\leq \sup \{|x^{**}(x^*)| : \|x^{**}\| \leq 1\} = \sup \{|J(x^*)(x^{**})| : \|x^{**}\| \leq 1\} \\ &\equiv \|Jx^*\| = \|x^*\|. \blacksquare\end{aligned}$$

Definition 21.2.14 When J maps X onto X'' , X is called reflexive.

It happens the L^p spaces are reflexive whenever $p > 1$. This is shown later.

21.3 Uniform Convexity of L^p

These terms refer roughly to how round the unit ball is. Here is the definition.

Definition 21.3.1 A Banach space is uniformly convex if whenever

$$\|x_n\|, \|y_n\| \leq 1$$

and $\|x_n + y_n\| \rightarrow 2$, it follows that $\|x_n - y_n\| \rightarrow 0$. More precisely, for every $\varepsilon > 0$, there is a $\delta > 0$ such that if $\|x + y\| > 2 - \delta$ for $\|x\|, \|y\|$ both no more than 1, then $\|x - y\| < \varepsilon$.

You can show that uniform convexity implies strict convexity. There are various other things which can also be shown. See the exercises for some of these. In this section, it will be shown that the L^p spaces are examples of uniformly convex spaces. This involves some inequalities known as Clarkson's inequalities. Before presenting these, here are the backwards Holder inequality and the backwards Minkowski inequality. Recall that in the Holder inequality, $\frac{p}{p-1} = q = p'$ and for $p > 1$,

$$\int_{\Omega} |f| |g| d\mu \leq \left(\int_{\Omega} |f|^p d\mu \right)^{1/p} \left(\int_{\Omega} |g|^{p/(p-1)} d\mu \right)^{(p-1)/p}$$

The idea in these inequalities is to consider the case that $p \in (0, 1)$. This inequality is easy to remember if you just take Holder's inequality and turn it around in the case that $0 < p < 1$.

Lemma 21.3.2 Let $0 < p < 1$ and let f, g be measurable functions. Also

$$\int_{\Omega} |g|^{p/(p-1)} d\mu < \infty, \int_{\Omega} |f|^p d\mu < \infty,$$

which implies that g is 0 only on a set of measure zero. Then the following backwards Holder inequality holds.

$$\int_{\Omega} |fg| d\mu \geq \left(\int_{\Omega} |f|^p d\mu \right)^{1/p} \left(\int_{\Omega} |g|^{p/(p-1)} d\mu \right)^{(p-1)/p}$$

Proof: If $\int |fg| d\mu = \infty$, there is nothing to prove. Hence assume this is finite. Then

$$\int |f|^p d\mu = \int |g|^{-p} |fg|^p d\mu$$

This makes sense because, due to the hypothesis on g it must be the case that g equals 0 only on a set of measure zero, since $p/(p-1) < 0$.

Then by the usual Holder inequality, one of the exponents being $1/p > 1$, the other being $1/(1-p)$ also larger than 1 with $p + (1-p) = 1$,

$$\begin{aligned} \int |f|^p d\mu &\leq \left(\int |fg| d\mu \right)^p \left(\int \left(\frac{1}{|g|^p} \right)^{1/(1-p)} d\mu \right)^{1-p} \\ &= \left(\int |fg| d\mu \right)^p \left(\int |g|^{p/(p-1)} d\mu \right)^{1-p} \end{aligned}$$

Now divide by $\left(\int |g|^{p/(p-1)} d\mu \right)^{1-p}$ and then take the p^{th} root. ■

Here is the backwards Minkowski inequality. It looks just like the ordinary Minkowski inequality except the inequality is turned around.

Corollary 21.3.3 *Let $0 < p < 1$ and suppose $\int |h|^p d\mu < \infty$ for $h = f, g$. Then*

$$\left(\int (|f| + |g|)^p d\mu \right)^{1/p} \geq \left(\int |f|^p d\mu \right)^{1/p} + \left(\int |g|^p d\mu \right)^{1/p}$$

Proof: If $\int (|f| + |g|)^p d\mu = 0$ then there is nothing to prove since this implies $|f| = |g| = 0$ a.e. so assume this is not zero.

$$\int (|f| + |g|)^p d\mu = \int (|f| + |g|)^{p-1} (|f| + |g|) d\mu$$

Since $p < 1$, $(|f| + |g|)^p \leq |f|^p + |g|^p$ and so

$$\int \left((|f| + |g|)^{p-1} \right)^{p/(p-1)} d\mu < \infty.$$

Hence the backward Holder inequality applies and it follows that

$$\begin{aligned} \int (|f| + |g|)^p d\mu &= \int (|f| + |g|)^{p-1} |f| d\mu + \int (|f| + |g|)^{p-1} |g| d\mu \\ &\geq \left(\int \left((|f| + |g|)^{p-1} \right)^{p/(p-1)} d\mu \right)^{(p-1)/p} \left[\left(\int |f|^p d\mu \right)^{1/p} + \left(\int |g|^p d\mu \right)^{1/p} \right] \\ &= \left(\int (|f| + |g|)^p d\mu \right)^{(p-1)/p} \left[\left(\int |f|^p d\mu \right)^{1/p} + \left(\int |g|^p d\mu \right)^{1/p} \right] \end{aligned}$$

and so, dividing gives the desired inequality. ■

Consider the “easy” Clarkson inequalities.

Lemma 21.3.4 For any $p \geq 2$ the following inequality holds for any $t \in [0, 1]$,

$$\left(\frac{1+t}{2}\right)^p + \left(\frac{1-t}{2}\right)^p \leq \frac{1}{2}(t^p + 1)$$

Proof: It is clear that, since $p \geq 2$, the inequality holds for $t = 0$ and $t = 1$. Thus it suffices to consider only $t \in (0, 1)$. Let $x = 1/t$. Then, dividing by t^p , the inequality holds if and only if

$$\left(\frac{x+1}{2}\right)^p + \left(\frac{x-1}{2}\right)^p \leq \frac{1}{2}(1+x^p)$$

for all $x \geq 1$. Let

$$f(x) = \frac{1}{2}(1+x^p) - \left(\left(\frac{x+1}{2}\right)^p + \left(\frac{x-1}{2}\right)^p\right)$$

Then $f(1) = 0$ and

$$f'(x) = \frac{p}{2}x^{p-1} - \left(\frac{p}{2}\left(\frac{x+1}{2}\right)^{p-1} + \frac{p}{2}\left(\frac{x-1}{2}\right)^{p-1}\right)$$

Since $p-1 \geq 1$, $g(x) = x^{p-1}$ is convex. Its graph is like a smile. Thus $\frac{1}{2}(g(x_1) + g(x_2)) \geq g\left(\frac{x_1+x_2}{2}\right)$ and so

$$f'(x) \geq \frac{p}{2}x^{p-1} - p\left(\frac{\frac{x+1}{2} + \frac{x-1}{2}}{2}\right)^{p-1} = \frac{p}{2}x^{p-1} - p\left(\frac{x}{2}\right)^{p-1} \geq 0$$

Hence $f(x) \geq 0$ for all $x \geq 1$. ■

Corollary 21.3.5 If $z, w \in \mathbb{C}$ and $p \geq 2$, then

$$\left|\frac{z+w}{2}\right|^p + \left|\frac{z-w}{2}\right|^p \leq \frac{1}{2}(|z|^p + |w|^p) \quad (21.6)$$

Proof: One of $|w|, |z|$ is larger. Say $|z| \geq |w|$. Then dividing both sides of the proposed inequality by $|z|^p$ it suffices to verify that for all complex t having $|t| \leq 1$,

$$\left|\frac{1+t}{2}\right|^p + \left|\frac{1-t}{2}\right|^p \leq \frac{1}{2}(|t|^p + 1) \quad (21.7)$$

Say $t = re^{i\theta}$ where $r \leq 1$. Then we need to estimate

$$\left|\frac{1+re^{i\theta}}{2}\right|^p + \left|\frac{1-re^{i\theta}}{2}\right|^p$$

It suffices to show that this is no larger than $\frac{1}{2}(r^p + 1)$. The function on the left in 21.7 equals

$$\begin{aligned} & \frac{1}{2^p} \left((1+r\cos\theta)^2 + r^2\sin^2(\theta) \right)^{p/2} + \left((1-r\cos\theta)^2 + r^2\sin^2(\theta) \right)^{p/2} \\ &= \frac{1}{2^p} (1+r^2+2r\cos\theta)^{p/2} + (1+r^2-2r\cos\theta)^{p/2}, \end{aligned} \quad (21.8)$$

I want to find the maximum value of this function of θ for $\theta \in [0, 2\pi]$. By calculus, this will be when the derivative is 0 or at an endpoint. The derivative with respect to θ is $\frac{1}{2^p}$ times

$$\begin{aligned} & \frac{p}{2} \left((1+r^2+2r\cos\theta)^{\frac{p-2}{2}} (-2r\sin\theta) + (2r\sin\theta) (1+r^2-2r\cos\theta)^{\frac{p-2}{2}} \right) \\ &= \frac{p}{2} (2r\sin\theta) \left((1+r^2-2r\cos\theta)^{\frac{p-2}{2}} - (1+r^2+2r\cos\theta)^{\frac{p-2}{2}} \right) \end{aligned}$$

This equals 0 when $\theta = 0, \pi, 2\pi$ or when $\theta = \frac{\pi}{2}, \frac{3\pi}{2}$. At the last two values, the value of the function in 21.8 is

$$\frac{1}{2^{p-1}} (1+r^2)^{p/2} \leq \frac{1}{2} (r^p + 1).$$

This follows from convexity of $y = x^{p/2}$ for $p \geq 2$. Here is why:

$$\frac{1}{2^{p-1}} (1+r^2)^{p/2} = \frac{2^{p/2}}{2^{p-1}} \left(\frac{1+r^2}{2} \right)^{p/2} \leq \frac{2^{p/2}}{2^{p-1}} \frac{1}{2} (1+r^p)$$

At 0 or π , the value of the function in 21.8 is

$$\left((1+r^2-2r)^{p/2} + (1+r^2+2r)^{p/2} \right) \frac{1}{2^p} = \left(\frac{1+r}{2} \right)^p + \left(\frac{1-r}{2} \right)^p$$

and from the above lemma, this is no larger than $\frac{1}{2} (r^p + 1)$. ■

With this corollary, here is the easy Clarkson inequality.

Theorem 21.3.6 *Let $p \geq 2$. Then*

$$\left\| \frac{f+g}{2} \right\|_{L^p}^p + \left\| \frac{f-g}{2} \right\|_{L^p}^p \leq \frac{1}{2} (\|f\|_{L^p}^p + \|g\|_{L^p}^p)$$

Proof: This follows right away from the above corollary.

$$\int_{\Omega} \left| \frac{f+g}{2} \right|^p d\mu + \int_{\Omega} \left| \frac{f-g}{2} \right|^p d\mu \leq \frac{1}{2} \int_{\Omega} (|f|^p + |g|^p) d\mu \quad \blacksquare$$

Now it remains to consider the hard Clarkson inequalities. These pertain to $p < 2$. First is the following elementary inequality.

Lemma 21.3.7 *For $1 < p < 2$, the following inequality holds for all $t \in [0, 1]$.*

$$\left(\frac{1+t}{2} \right)^q + \left(\frac{1-t}{2} \right)^q \leq \left(\frac{1}{2} + \frac{1}{2} t^p \right)^{q/p}$$

where here $1/p + 1/q = 1$ so $q > 2$.

Proof: First note that if $t = 0$ or 1 , the inequality holds. Next observe that the map $s \rightarrow \frac{1-s}{1+s}$ maps $(0, 1)$ onto $(0, 1)$. Replace t with $(1-s)/(1+s)$. Then the desired inequality is equivalent to the following for $s \in (0, 1)$.

$$\left(\frac{1}{s+1} \right)^q + \left(\frac{s}{s+1} \right)^q \leq \left(\frac{1}{2} + \frac{1}{2} \left(\frac{1-s}{s+1} \right)^p \right)^{q/p}$$

Multiplying both sides by $(1+s)^q$, this inequality is equivalent to showing that for all $s \in (0, 1)$,

$$\begin{aligned} 1 + s^q &\leq ((1+s)^p)^{q/p} \left(\frac{1}{2} + \frac{1}{2} \left(\frac{1-s}{s+1} \right)^p \right)^{q/p} \\ &= \left(\frac{1}{2} \right)^{q/p} ((1+s)^p + (1-s)^p)^{q/p} \end{aligned}$$

This is the same as establishing

$$\frac{1}{2} ((1+s)^p + (1-s)^p) - (1+s^q)^{p-1} \geq 0 \quad (21.9)$$

where $p-1 = p/q$ due to the definition of q above. Note how this has reduced to an expression in which exponents are p or $p-1$ rather than q . We know $p \in (1, 2)$ whereas, q is something larger than 2.

$$\binom{p}{l} \equiv \frac{p(p-1)\cdots(p-l+1)}{l!}, \quad l \geq 1$$

and $\binom{p}{0} \equiv 1$. What is the sign of $\binom{p}{l}$? Recall that $1 < p < 2$ so the sign is positive if $l = 0, l = 1, l = 2$. What about $l = 3$? $\binom{p}{3} = \frac{p(p-1)(p-2)}{3!}$ so this is negative. Then $\binom{p}{4}$ is positive. Thus these alternate between positive and negative with $\binom{p}{2k} > 0$ for all k . What about $\binom{p-1}{k}$? When $k = 0$ it is positive. When $k = 1$ it is also positive. When $k = 2$ it equals $\frac{(p-1)(p-2)}{2!} < 0$. Then when $k = 3$, $\binom{p-1}{3} > 0$. Thus $\binom{p-1}{k}$ is positive when k is odd and is negative when k is even.

Now return to 21.9. The left side equals

$$\frac{1}{2} \left(\sum_{k=0}^{\infty} \binom{p}{k} s^k + \sum_{k=0}^{\infty} \binom{p}{k} (-s)^k \right) - \sum_{k=0}^{\infty} \binom{p-1}{k} s^{qk}.$$

The first term equals 0. Then this reduces to

$$\sum_{k=1}^{\infty} \binom{p}{2k} s^{2k} - \binom{p-1}{2k} s^{q2k} - \binom{p-1}{2k-1} s^{q(2k-1)}$$

From the above observation about the binomial coefficients, the above is larger than

$$\sum_{k=1}^{\infty} \binom{p}{2k} s^{2k} - \binom{p-1}{2k-1} s^{q(2k-1)}$$

It remains to show the k^{th} term in the above sum is nonnegative. Now $q(2k-1) > 2k$ for all $k \geq 1$ because $q > 2$. Then since $0 < s < 1$

$$\binom{p}{2k} s^{2k} - \binom{p-1}{2k-1} s^{q(2k-1)} \geq s^{2k} \left(\binom{p}{2k} - \binom{p-1}{2k-1} \right)$$

However, this is nonnegative because it equals

$$\begin{aligned}
 & s^{2k} \left(\frac{p(p-1)\cdots(p-2k+1)}{(2k)!} - \overbrace{\frac{(p-1)(p-2)\cdots(p-2k+1)}{(2k-1)!}}^{>0} \right) \\
 & \geq s^{2k} \left(\frac{p(p-1)\cdots(p-2k+1)}{(2k)!} - \frac{(p-1)(p-2)\cdots(p-2k+1)}{(2k)!} \right) \\
 & = s^{2k} \frac{(p-1)(p-2)\cdots(p-2k+1)}{(2k)!} (p-1) > 0. \blacksquare
 \end{aligned}$$

Corollary 21.3.8 *Let $z, w \in \mathbb{C}$. Then for $p \in (1, 2)$,*

$$\left| \frac{z+w}{2} \right|^q + \left| \frac{z-w}{2} \right|^q \leq \left(\frac{1}{2} |z|^p + \frac{1}{2} |w|^p \right)^{q/p}$$

Proof: One of $|w|, |z|$ is larger. Say $|w| \geq |z|$. Then dividing by $|w|^q$, for $t = z/w$, showing the above inequality is equivalent to showing that for all $t \in \mathbb{C}$, $|t| \leq 1$,

$$\left| \frac{t+1}{2} \right|^q + \left| \frac{1-t}{2} \right|^q \leq \left(\frac{1}{2} |t|^p + \frac{1}{2} \right)^{q/p}$$

Now $q > 2$ and so by the same argument given in proving Corollary 21.3.5, for $t = re^{i\theta}$, the left side of the above inequality is maximized when $\theta = 0$. Hence, from Lemma 21.3.7,

$$\begin{aligned}
 \left| \frac{t+1}{2} \right|^q + \left| \frac{1-t}{2} \right|^q & \leq \left| \frac{|t|+1}{2} \right|^q + \left| \frac{1-|t|}{2} \right|^q \\
 & \leq \left(\frac{1}{2} |t|^p + \frac{1}{2} \right)^{q/p}. \blacksquare
 \end{aligned}$$

From this the hard Clarkson inequality follows. The two Clarkson inequalities are summarized in the following theorem.

Theorem 21.3.9 *Let $2 \leq p$. Then*

$$\left\| \frac{f+g}{2} \right\|_{L^p}^p + \left\| \frac{f-g}{2} \right\|_{L^p}^p \leq \frac{1}{2} (\|f\|_{L^p}^p + \|g\|_{L^p}^p)$$

Let $1 < p < 2$. Then for $1/p + 1/q = 1$,

$$\left\| \frac{f+g}{2} \right\|_{L^p}^q + \left\| \frac{f-g}{2} \right\|_{L^p}^q \leq \left(\frac{1}{2} \|f\|_{L^p}^p + \frac{1}{2} \|g\|_{L^p}^p \right)^{q/p}$$

Proof: The first was established above. Consider the second.

$$\begin{aligned}
 & \left\| \frac{f+g}{2} \right\|_{L^p}^q + \left\| \frac{f-g}{2} \right\|_{L^p}^q = \\
 & \left(\int_{\Omega} \left| \frac{f+g}{2} \right|^p d\mu \right)^{q/p} + \left(\int_{\Omega} \left| \frac{f-g}{2} \right|^p d\mu \right)^{q/p}
 \end{aligned}$$

$$= \left(\int_{\Omega} \left(\left| \frac{f+g}{2} \right|^q \right)^{p/q} d\mu \right)^{q/p} + \left(\int_{\Omega} \left(\left| \frac{f-g}{2} \right|^q \right)^{p/q} d\mu \right)^{q/p}$$

Now $p/q < 1$ and so the backwards Minkowski inequality applies. Thus

$$\leq \left(\int_{\Omega} \left(\left| \frac{f+g}{2} \right|^q + \left| \frac{f-g}{2} \right|^q \right)^{p/q} d\mu \right)^{q/p}$$

From Corollary 21.3.8,

$$\begin{aligned} &\leq \left(\int_{\Omega} \left(\left(\frac{1}{2} |f|^p + \frac{1}{2} |g|^p \right)^{q/p} \right)^{p/q} d\mu \right)^{q/p} \\ &= \left(\int_{\Omega} \left(\frac{1}{2} |f|^p + \frac{1}{2} |g|^p \right) d\mu \right)^{q/p} = \left(\frac{1}{2} \|f\|_{L^p}^p + \frac{1}{2} \|g\|_{L^p}^p \right)^{q/p} \blacksquare \end{aligned}$$

Now with these Clarkson inequalities, it is not hard to show that all the L^p spaces are uniformly convex.

Theorem 21.3.10 *The L^p spaces are uniformly convex.*

Proof: First suppose $p \geq 2$. Suppose $\|f_n\|_{L^p}, \|g_n\|_{L^p} \leq 1$ and $\left\| \frac{f_n+g_n}{2} \right\|_{L^p} \rightarrow 1$. Then from the first Clarkson inequality,

$$\left\| \frac{f_n+g_n}{2} \right\|_{L^p}^p + \left\| \frac{f_n-g_n}{2} \right\|_{L^p}^p \leq \frac{1}{2} (\|f_n\|_{L^p}^p + \|g_n\|_{L^p}^p) \leq 1$$

and so $\|f_n - g_n\|_{L^p} \rightarrow 0$.

Next suppose $1 < p < 2$ and $\left\| \frac{f_n+g_n}{2} \right\|_{L^p} \rightarrow 1$. Then from the second Clarkson inequality

$$\left\| \frac{f_n+g_n}{2} \right\|_{L^p}^q + \left\| \frac{f_n-g_n}{2} \right\|_{L^p}^q \leq \left(\frac{1}{2} \|f_n\|_{L^p}^p + \frac{1}{2} \|g_n\|_{L^p}^p \right)^{q/p} \leq 1$$

which shows that $\|f_n - g_n\|_{L^p} \rightarrow 0$. \blacksquare

21.4 Closed Subspaces

Theorem 21.4.1 *Let X be a Banach space and let $V = \text{span}(x_1, \dots, x_n)$. Then V is a closed subspace of X .*

Proof: Without loss of generality, it can be assumed $\{x_1, \dots, x_n\}$ is linearly independent. Otherwise, delete those vectors which are in the span of the others till a linearly independent set is obtained. Let

$$x = \lim_{p \rightarrow \infty} \sum_{k=1}^n c_k^p x_k \in \overline{V}. \quad (21.10)$$

First suppose $\mathbf{c}^p \equiv (c_1^p, \dots, c_n^p)$ is not bounded in \mathbb{F}^n . Then $\mathbf{d}^p \equiv \mathbf{c}^p / \|\mathbf{c}^p\|_{\mathbb{F}^n}$ is a unit vector in \mathbb{F}^n and so there exists a subsequence, still denoted by \mathbf{d}^p which converges to \mathbf{d} where $\|\mathbf{d}\| = 1$. Then

$$\mathbf{0} = \lim_{p \rightarrow \infty} \frac{x}{\|\mathbf{c}^p\|} = \lim_{p \rightarrow \infty} \sum_{k=1}^n d_k^p x_k = \sum_{k=1}^n d_k x_k$$

where $\sum_k |d_k|^2 = 1$ in contradiction to the linear independence of the $\{x_1, \dots, x_n\}$. Hence it must be the case that \mathbf{c}^p is bounded in \mathbb{F}^n . Then taking a subsequence, still denoted as p , it can be assumed $\mathbf{c}^p \rightarrow \mathbf{c}$ and then in 21.10 it follows $x = \sum_{k=1}^n c_k x_k \in \text{span}(x_1, \dots, x_n)$. ■

Proposition 21.4.2 *Let E be a separable Banach space. Then there exists an increasing sequence of subspaces, $\{F_n\}$ such that $\dim(F_{n+1}) - \dim(F_n) \leq 1$ and equals 1 for all n if the dimension of E is infinite. Also $\cup_{n=1}^{\infty} F_n$ is dense in E . In the case where E is infinite dimensional, $F_n = \text{span}(e_1, \dots, e_n)$ where for each n*

$$\text{dist}(e_{n+1}, F_n) \geq \frac{1}{2} \quad (21.11)$$

and defining,

$$G_k \equiv \text{span}(\{e_j : j \neq k\})$$

$$\text{dist}(e_k, G_k) \geq \frac{1}{4}. \quad (21.12)$$

Proof: Since E is separable, so is $\partial B(0, 1)$, the boundary of the unit ball thanks to Corollary 3.4.3. Let $\{w_k\}_{k=1}^{\infty}$ be a countable dense subset of $\partial B(0, 1)$.

Let $e_1 = w_1$. Let $F_1 = \mathbb{F}e_1$. Suppose F_n has been obtained and equals the following: $\text{span}(e_1, \dots, e_n)$ where $\{e_1, \dots, e_n\}$ is independent, $\|e_k\| = 1$, and

$$\text{dist}(e_n, \text{span}(e_1, \dots, e_{n-1})) \geq \frac{1}{2}.$$

For each n , F_n is closed by Theorem 21.4.1.

If F_n contains $\{w_k\}_{k=1}^{\infty}$, let $F_m = F_n$ for all $m > n$. Otherwise, pick $w \in \{w_k\}$ to be the point of $\{w_k\}_{k=1}^{\infty}$ having the smallest subscript which is not contained in F_n . Then w is at a positive distance λ from F_n because F_n is closed. Therefore, there exists $y \in F_n$ such that $\lambda \leq \|y - w\| \leq 2\lambda$. Let $e_{n+1} = \frac{w-y}{\|w-y\|}$. It follows

$$w = \|w-y\| e_{n+1} + y \in \text{span}(e_1, \dots, e_{n+1}) \equiv F_{n+1}$$

Then if $x \in \text{span}(e_1, \dots, e_n)$,

$$\begin{aligned} \|e_{n+1} - x\| &= \left\| \frac{w-y}{\|w-y\|} - x \right\| = \left\| \frac{w-y}{\|w-y\|} - \frac{\|w-y\| x}{\|w-y\|} \right\| \\ &\geq \frac{1}{2\lambda} \|w-y - \|w-y\| x\| \geq \frac{\lambda}{2\lambda} = \frac{1}{2}. \end{aligned}$$

This has shown the existence of an increasing sequence of subspaces, $\{F_n\}$ as described above. It remains to show the union of these subspaces is dense. First note that the union of these subspaces must contain the $\{w_k\}_{k=1}^{\infty}$ because if w_m is missing, then it would contradict

the construction at the m^{th} step. That one should have been chosen. However, $\{w_k\}_{k=1}^{\infty}$ is dense in $\partial B(0, 1)$. If $x \in E$ and $x \neq 0$, then $\frac{x}{\|x\|} \in \partial B(0, 1)$ then there exists

$$w_m \in \{w_k\}_{k=1}^{\infty} \subseteq \bigcup_{n=1}^{\infty} F_n$$

such that $\left\|w_m - \frac{x}{\|x\|}\right\| < \frac{\varepsilon}{\|x\|}$. But then $\| \|x\| w_m - x \| < \varepsilon$. and so $\|x\| w_m$ is a point of $\bigcup_{n=1}^{\infty} F_n$ which is within ε of x . This proves $\bigcup_{n=1}^{\infty} F_n$ is dense as desired. 21.11 follows from the construction. It remains to verify 21.12.

Let $y \in G_k$. Thus for some $n, y = \sum_{j=1}^{k-1} c_j e_j + \sum_{j=k+1}^n c_j e_j$ and I need to show $\|y - e_k\| \geq 1/4$. Without loss of generality, $c_n \neq 0$ and $n > k$. Suppose 21.12 does not hold for some such y so that

$$\left\| e_k - \left(\sum_{j=1}^{k-1} c_j e_j + \sum_{j=k+1}^n c_j e_j \right) \right\| < \frac{1}{4}. \quad (21.13)$$

Then from the construction,

$$\frac{1}{4} > |c_n| \left\| e_k - \left(\sum_{j=1}^{k-1} (c_j/c_n) e_j + \sum_{j=k+1}^{n-1} (c_j/c_n) e_j + e_n \right) \right\| \geq |c_n| \frac{1}{2}$$

and so $|c_n| < 1/2$. Consider the left side of 21.13. By the construction

$$\begin{aligned} \frac{1}{4} &> \left\| \overbrace{c_n(e_k - e_n) + (1 - c_n)e_k}^{e_k - c_n e_n} - \left(\sum_{j=1}^{k-1} c_j e_j + \sum_{j=k+1}^{n-1} c_j e_j \right) \right\| \\ &\geq |1 - c_n| - |c_n| \left\| (e_k - e_n) - \left(\sum_{j=1}^{k-1} (c_j/c_n) e_j + \sum_{j=k+1}^{n-1} (c_j/c_n) e_j \right) \right\| \\ &\geq |1 - c_n| - |c_n| \frac{1}{2} \geq 1 - \frac{3}{2} |c_n| > 1 - \frac{3}{2} \frac{1}{2} = \frac{1}{4}, \end{aligned}$$

a contradiction. This proves the desired estimate. ■

Definition 21.4.3 A Banach space X has a Schauder basis $\{e_k\}_{k=1}^{\infty}$ if for every $x \in X$, there are unique scalars c_k such that $x = \sum_{k=1}^{\infty} c_k x_k$. This is different than a basis because you allow countable sums. For example, you might consider Fourier series.

21.5 Weak And Weak * Topologies

Proposition 21.4.2 shows that in infinite dimensional space, closed and bounded will not be compact. However, in applications one would like to be able to get convergence of subsequences. This involves asking for less than norm convergence and the concept of weak topologies.

21.5.1 Basic Definitions

Let X be a Banach space and let X' be its dual space.¹ For A' a **finite** subset of X' , denote by $\rho_{A'}$ the function defined on X

$$\rho_{A'}(x) \equiv \max_{x^* \in A'} |x^*(x)| \quad (21.14)$$

¹ Actually, all this works in much more general settings than this.

and also let $B_{A'}(x, r)$ be defined by

$$B_{A'}(x, r) \equiv \{y \in X : \rho_{A'}(y - x) < r\} \quad (21.15)$$

Then certain things are obvious. First of all, if $a \in \mathbb{F}$ and $x, y \in X$,

$$\begin{aligned} \rho_{A'}(x + y) &\leq \rho_{A'}(x) + \rho_{A'}(y), \\ \rho_{A'}(ax) &= |a| \rho_{A'}(x). \end{aligned}$$

Similarly, letting A be a finite subset of X , denote by ρ_A the function defined on X'

$$\rho_A(x^*) \equiv \max_{x \in A} |x^*(x)| \quad (21.16)$$

and let $B_A(x^*, r)$ be defined by

$$B_A(x^*, r) \equiv \{y^* \in X' : \rho_A(y^* - x^*) < r\}. \quad (21.17)$$

It is also clear that

$$\begin{aligned} \rho_A(x^* + y^*) &\leq \rho_A(x^*) + \rho_A(y^*), \\ \rho_A(ax^*) &= |a| \rho_A(x^*). \end{aligned}$$

Lemma 21.5.1 *The sets $B_{A'}(x, r)$ where A' is a finite subset of X' and $x \in X$ form a basis for a topology on X known as the weak topology. The sets $B_A(x^*, r)$ where A is a finite subset of X and $x^* \in X'$ form a basis for a topology on X' known as the weak * topology.*

Proof: The two assertions are very similar. I will verify the one for the weak topology. The union of these sets, $B_{A'}(x, r)$ for $x \in X$ and $r > 0$ is all of X . Now suppose z is contained in the intersection of two of these sets. Say

$$z \in B_{A'}(x, r) \cap B_{A'_1}(x_1, r_1)$$

Then let $C' = A' \cup A'_1$ and let

$$0 < \delta \leq \min(r - \rho_{A'}(z - x), r_1 - \rho_{A'_1}(z - x_1)).$$

Consider $y \in B_{C'}(z, \delta)$. Then

$$r - \rho_{A'}(z - x) \geq \delta > \rho_{C'}(y - z) \geq \rho_{A'}(y - z)$$

and so

$$r > \rho_{A'}(y - z) + \rho_{A'}(z - x) \geq \rho_{A'}(y - x)$$

which shows $y \in B_{A'}(x, r)$. Similar reasoning shows $y \in B_{A'_1}(x_1, r_1)$ and so

$$B_{C'}(z, \delta) \subseteq B_{A'}(x, r) \cap B_{A'_1}(x_1, r_1).$$

Therefore, these sets are a basis for a topology known as the weak topology which consists of the union of all sets of the form $B_{A'}(x, r)$. ■

21.5.2 Banach Alaoglu Theorem

Why does anyone care about these topologies? The short answer is that in the weak $*$ topology, the closed unit ball in X' is compact. This is not true in the norm topology thanks to Proposition 21.4.2. This wonderful result is the Banach Alaoglu theorem. First recall the notion of the product topology, and the Tychonoff theorem, Theorem 19.3.2 on Page 509 which are stated here for convenience.

Definition 21.5.2 Let I be a set and suppose for each $i \in I$, (X_i, τ_i) is a nonempty topological space. The Cartesian product of the X_i , denoted by $\prod_{i \in I} X_i$, consists of the set of all choice functions defined on I which select a single element of each X_i . Thus $f \in \prod_{i \in I} X_i$ means for every $i \in I$, $f(i) \in X_i$. The axiom of choice says $\prod_{i \in I} X_i$ is nonempty. Next is a description of a subbasis for a topology. Let $P_j(A) = \prod_{i \in I} B_i$ where $B_i = X_i$ if $i \neq j$ and $B_j = A$. A subbasis for a topology on the product space consists of all sets $P_j(A)$ where $A \in \tau_j$. (These sets have an open set from the topology of X_j in the j^{th} slot and the whole space in the other slots.) Thus a basis consists of finite intersections of these sets. Note that the intersection of two of these basic sets is another basic set and their union yields $\prod_{i \in I} X_i$. Therefore, they satisfy the condition needed for a collection of sets to serve as a basis for a topology. This topology is called the product topology and is denoted by $\prod \tau_i$.

Theorem 21.5.3 If (X_i, τ_i) is compact, then so is $(\prod_{i \in I} X_i, \prod \tau_i)$.

The Banach Alaoglu theorem is as follows.

Theorem 21.5.4 Let B' be the closed unit ball in X' . Then B' is compact in the weak $*$ topology.

Proof: By the Tychonoff theorem, Theorem 21.5.3, $P \equiv \prod_{x \in X} \overline{B(0, \|x\|)}$ is compact in the product topology where the topology on $\overline{B(0, \|x\|)}$ is the usual topology of \mathbb{F} . Recall P is the set of functions which map a point $x \in X$ to a point in $\overline{B(0, \|x\|)}$. Therefore, $B' \subseteq P$. Also the basic open sets in the weak $*$ topology on B' are obtained as the intersection of basic open sets in the product topology of P to B' and so it suffices to show B' is a closed subset of P . Suppose then that $f \in P \setminus B'$. Since $|f(x)| \leq \|x\|$ for each x , it follows f cannot be linear. There are two ways this can happen. One way is that for some x, y

$$f(x+y) \neq f(x) + f(y) \quad (21.18)$$

for some $x, y \in X$ and the other is that $f(\lambda x) \neq \lambda f(x)$ for some λ, x . Consider the first. If g is close enough to f at the three points, $x+y, x$, and y , 21.18 will hold for g in place of f . In other words there is a basic open set containing f , such that for all g in this basic open set, $g \notin B'$. A similar consideration applies in case $f(\lambda x) \neq \lambda f(x)$ for some scalar λ and x . Since $P \setminus B'$ is open, it follows B' is a closed subset of P and is therefore, compact. ■

Note that if the canonical map $J: X \rightarrow X''$ discussed earlier given by $Jx(x^*) \equiv x^*(x)$ is onto, then we could conclude that B the closed unit ball in X is weakly compact because J would be a homeomorphism of B'' and B . Thus reflexive spaces are important in these considerations.

Sometimes one can consider the weak $*$ topology as a metric space. You can do it for K when K is weak $*$ compact and X is separable. Note that it was just shown that the closed ball is weak $*$ compact.

Theorem 21.5.5 *If $K \subseteq X'$ is compact in the weak * topology and X is separable in the weak topology then there exists a metric d , on K such that if τ_d is the topology on K induced by d and if τ is the topology on K induced by the weak * topology of X' , then $\tau = \tau_d$. Thus one can consider K with the weak * topology as a metric space.*

Proof: Let $D = \{x_n\}$ be the dense countable subset in X . The metric is

$$d(f, g) \equiv \sum_{n=1}^{\infty} 2^{-n} \frac{\rho_{x_n}(f - g)}{1 + \rho_{x_n}(f - g)}$$

where $\rho_{x_n}(f) = |f(x_n)|$. Clearly $d(f, g) = d(g, f) \geq 0$. If $d(f, g) = 0$, then this requires $f(x_n) = g(x_n)$ for all $x_n \in D$. Is it the case that $f = g$? Does $f(x) = g(x)$ for all x , not just for the x_n ?

Letting x be given, $B_{\{f, g\}}(x, r)$ contains some $x_n \in D$. Hence

$$\max\{|f(x_n) - f(x)|, |g(x_n) - g(x)|\} < r$$

and $f(x_n) = g(x_n)$. It follows that $|f(x) - g(x)| \leq$

$$|f(x) - f(x_n)| + \overbrace{|f(x_n) - g(x_n)|}^{=0} + |g(x_n) - g(x)| < 2r.$$

Since r is arbitrary, this implies $f(x) = g(x)$.

It is routine to verify the triangle inequality from the easy to establish inequality,

$$\frac{x}{1+x} + \frac{y}{1+y} \geq \frac{x+y}{1+x+y},$$

valid whenever $x, y \geq 0$. Therefore this is a metric.

Thus there are two topological spaces, (K, τ) and (K, d) , the first being K with the weak * topology and the second being K with the topology from this metric. Suppose $B(f, r)$ is an open ball with respect to the metric space topology. I claim that $B(f, r)$ is open in the weak * topology τ . To do this, let $d(f, g) < r$. Is there a finite set $A \subseteq X$ such that $B_A(g, \delta) \subseteq B(f, r)$? Let $A_n \equiv \{x_1, \dots, x_n\}$ and pick n large enough that $\sum_{k=n}^{\infty} 2^{-k} < \frac{r-d(f, g)}{2} \equiv \delta$. Then if $h \in B_{A_n}(g, \delta)$, it follows that

$$\begin{aligned} d(f, h) &\leq d(f, g) + d(g, h) < d(f, g) + \sum_{k=1}^{n-1} 2^{-k} |g(x_k) - h(x_k)| + \frac{r-d(f, g)}{2} \\ &< d(f, g) + \sum_{k=1}^{\infty} \delta 2^{-k} + \delta = d(f, g) + 2\delta = d(f, g) + (r-d(f, g)) = r \end{aligned}$$

Thus $B_A(g, \delta) \subseteq B(f, r)$ and so $B(f, r)$ is the union of weak * open sets and is therefore, weakly open. It follows that $\tau_d \subseteq \tau$. Thus it is clear that if i is the identity map, $i: (K, \tau) \rightarrow (K, d)$, then i is continuous.

Now suppose $U \in \tau$. Is U in τ_d ? Since K is compact with respect to τ , it follows from the above that K is compact with respect to $\tau_d \subseteq \tau$. Hence $K \setminus U$ is compact with respect to τ_d and so it is closed with respect to τ_d . Thus U is open with respect to τ_d . The identity map $i: (K, d) \rightarrow (K, \tau)$ is continuous. ■

Note that the above proof is about the elements of X' continuous with respect to the weak topology on X and D is dense with respect to this **weak topology**.

The fact that this set with the weak $*$ topology can be considered a metric space is very significant because if a point is a limit point in a metric space, one can extract a convergent sequence. Also this has shown that the closed unit ball in X' can be considered a metric space provided X is separable. Therefore, it is sequentially compact from Theorem 3.5.8.

Note that if a Banach space is separable, then it is weakly separable. In fact, the countable dense set D with respect to the norm is also dense with respect to the weak topology. To see this, suppose $A = \{x_1^*, \dots, x_m^*\}$ is a finite subset of X' . Consider $B_A(x, r)$. Is there a point of D in $B_A(x, r)$? Choose $x_n \in D$ close enough to x that

$$\max \{|x_i^*(x) - x_i^*(x_n)|, i = 1, \dots, m\} < r$$

This can be done because $|x_i^*(x) - x_i^*(x_n)| \leq \|x_i^*\| \|x - x_n\|$ and this can be made small by taking

$$\|x - x_n\| < \min \left\{ \frac{r}{1 + \|x_i^*\|}, i = 1, \dots, m \right\}.$$

Corollary 21.5.6 *If X is weakly separable and $K \subseteq X'$ is compact in the weak $*$ topology, then K is sequentially compact. That is, if $\{f_n\}_{n=1}^\infty \subseteq K$, then there exists a subsequence f_{n_k} and $f \in K$ such that for all $x \in X$, $\lim_{k \rightarrow \infty} f_{n_k}(x) = f(x)$.*

Proof: By Theorem 21.5.5, K is a metric space for the metric described there and it is compact. Therefore by the characterization of compact metric spaces, Proposition 3.5.8 on Page 78, K is sequentially compact. This proves the corollary. ■

21.5.3 Eberlein Smulian Theorem

Next consider the weak topology. The most interesting results have to do with a reflexive Banach space. The following lemma ties together the weak and weak $*$ topologies in the case of a reflexive Banach space. It shows that a reflexive Banach space is actually a dual space.

Definition 21.5.7 *For X a Banach space, define $J : X \rightarrow X''$ by: For $x^* \in X'$, $Jx(x^*) \equiv x^*(x)$.*

For the properties of J see Theorem 21.2.13.

Lemma 21.5.8 *Let $J : X \rightarrow X''$ be the James map $Jx(x^*) \equiv x^*(x)$ and let X be reflexive so that J is onto. Then $J : (X, \text{weak topology}) \rightarrow (X'', \text{weak } * \text{ topology})$ is a homeomorphism. This means J is one to one, onto, and both J and J^{-1} are continuous.*

Proof: Let $x^* \in X'$ and let $B_{x^*}(x, r) \equiv \{y : |x^*(x) - x^*(y)| < r\}$. Thus $B_{x^*}(x, r)$ is a subbasic set for the weak topology on X . I claim that $JB_{x^*}(x, r) = B_{x^*}(Jx, r)$. where $B_{x^*}(Jx, r)$ is a subbasic set for the weak $*$ topology on X'' . If $y \in B_{x^*}(x, r)$, then $\|Jy - Jx\| = \|x - y\| < r$ and so $JB_{x^*}(x, r) \subseteq B_{x^*}(Jx, r)$. Now if $x^{**} \in B_{x^*}(Jx, r)$, then since J is reflexive, there exists $y \in X$ such that $Jy = x^{**}$ and so $\|y - x\| = \|Jy - Jx\| < r$ showing that $JB_{x^*}(x, r) = B_{x^*}(Jx, r)$. A typical subbasic set in the weak $*$ topology is of the form $B_{x^*}(Jx, r)$. Thus J maps the subbasic sets of the weak topology to the subbasic sets of the weak $*$ topology of X'' . Therefore, J is a homeomorphism as claimed. ■

The following is an easy corollary.

Corollary 21.5.9 *If X is a reflexive Banach space, then the closed unit ball is weakly compact.*

Proof: Let B be the closed unit ball. Then $B = J^{-1}(B^{**})$ where B^{**} is the unit ball in X'' which is compact in the weak * topology. Therefore B is weakly compact because J^{-1} is continuous. ■

Corollary 21.5.10 *If X is a reflexive Banach space, and X' is weak * separable, then the closed unit ball is weakly sequentially compact.*

Proof: This follows from Corollary 21.5.6 because X can be considered as the dual space of X' according to the rule $x(x^*) \equiv x^*(x)$. This definition gives x as continuous and linear. It is clearly linear and is continuous because $|x(x^*)| = |x^*(x)| \leq \|x^*\| \|x\|$. If $f \in X''$, then $f = Jx$ and so $f(x^*) = Jx(x^*) \equiv x^*(x) \equiv x(x^*)$. Since X' is weakly separable, it follows from the above corollary. To reiterate the reasoning, B is the unit ball in X'' and so it is weak * compact and is also a metric space so it is weakly sequentially compact also. ■

In fact if $K \subseteq X$ is weakly compact and X is reflexive with X' separable, then K is sequentially weakly compact.

Corollary 21.5.11 *Let X be a reflexive Banach space. If $K \subseteq X$ is compact in the weak topology and X' is separable in the weak * topology, then there exists a metric d , on K such that if τ_d is the topology on K induced by d and if τ is the topology on K induced by the weak topology of X , then $\tau = \tau_d$. Thus one can consider K with the weak topology as a metric space. Thus K is weakly sequentially compact.*

Proof: This follows from Theorem 21.5.5 and Lemma 21.5.8. Lemma 21.5.8 implies $J(K)$ is compact in X'' . Then since X' is separable in the weak * topology, X is separable in the weak topology and so there is a metric, d'' on $J(K)$ which delivers the weak * topology on $J(K)$. Let $d(x, y) \equiv d''(Jx, Jy)$. Then

$$(K, \tau_d) \xrightarrow{J} (J(K), \tau_{d''}) \xrightarrow{id} (J(K), \tau_{\text{weak}*}) \xrightarrow{J^{-1}} (K, \tau_{\text{weak}})$$

and all the maps are homeomorphisms. ■

Recall Lemma 21.2.9.

Lemma 21.5.12 *Let Y be a closed subspace of a Banach space X and let $y \in X \setminus Y$. Then there exists $x^* \in X'$ such that $x^*(Y) = 0$ but $x^*(y) \neq 0$.*

Next is the Eberlein Smulian theorem which states that a Banach space is reflexive if and only if the closed unit ball is weakly sequentially compact. Actually, only half the theorem is proved here, the more useful only if part. The book by Yoshida [60] has the complete theorem discussed. First here is an interesting lemma for its own sake.

Lemma 21.5.13 *A closed subspace of a reflexive Banach space is reflexive.*

Proof: Let Y be the closed subspace of the reflexive space, X . Consider the following diagram

$$\begin{array}{ccc} Y'' & \xrightarrow{i^{**} 1^{-1}} & X'' \\ Y' & \xleftarrow{i^* \text{ onto}} & X' \\ Y & \xrightarrow{i} & X \end{array}$$

This diagram follows from Theorem 21.2.10 on Page 544, the theorem on adjoints. Now let $y^{**} \in Y''$. Is $y^{**} = J_Y y$ for some $y \in Y$? Since X is reflexive, $i^{**} y^{**} = J_X(y)$ for some y . I want to show that $y \in Y$. If it is not in Y then since Y is closed, there exists $x^* \in X'$ such that $x^*(y) \neq 0$ but $x^*(Y) = 0$. Then $i^* x^* = 0$. Hence

$$0 = y^{**}(i^* x^*) = i^{**} y^{**}(x^*) = J(y)(x^*) = x^*(y) \neq 0,$$

a contradiction. Hence $y \in Y$. Letting J_Y denote the James map from Y to Y'' and $x^* \in X'$,

$$\begin{aligned} y^{**}(i^* x^*) &= i^{**} y^{**}(x^*) = J_X(y)(x^*) \\ &= x^*(y) = x^*(iy) = i^* x^*(y) = J_Y(y)(i^* x^*) \end{aligned}$$

Since i^* is onto, this shows $y^{**} = J_Y(y)$. ■

Theorem 21.5.14 (Eberlein Smulian) *The closed unit ball in a reflexive Banach space X , is weakly sequentially compact. By this is meant that if $\{x_n\}$ is contained in the closed unit ball, there exists a subsequence, $\{x_{n_k}\}$ and $x \in X$ such that for all $x^* \in X'$, it follows that $x^*(x_{n_k}) \rightarrow x^*(x)$.*

Proof: Let $\{x_n\} \subseteq B \equiv \overline{B(0, 1)}$. Let Y be the closure of the linear span of $\{x_n\}$. Thus Y is a separable. It is reflexive because it is a closed subspace of a reflexive space so the above lemma applies. By the Banach Alaoglu theorem, the closed unit ball B_Y^* in Y' is weak * compact. Also by Theorem 21.5.5, B_Y^* is a metric space with a suitable metric. The following diagram illustrates the rest of the argument.

$$\begin{array}{ccc} B^{**} & Y'' & \xrightarrow{i^{**} 1-1} X'' \\ B_Y^* & Y' \text{ weak * separable} & \xleftarrow{i^* \text{ onto}} X' \\ B_Y & Y \text{ separable} & \xrightarrow{i} X \end{array}$$

Thus B_Y^* is complete and totally bounded with respect to this metric and it follows that B_Y^* with the weak * topology is separable. This implies Y' is also separable in the weak * topology. To see this, let $\{y_n^*\} \equiv D$ be a weak * dense set in B_Y^* and let $y^* \in Y'$. Let p be a large enough positive rational number that $y^*/p \in B^*$. Then if A is any finite set from Y , there exists $y_n^* \in D$ such that $\rho_A(y^*/p - y_n^*) < \frac{\varepsilon}{p}$. It follows $py_n^* \in B_A(y^*, \varepsilon)$ showing that rational multiples of D are weak * dense in Y' . Letting $B_Y = B \cap Y$, this B_Y is the closed unit ball in Y and Y' is weak * separable. Therefore, by Corollary 21.5.10, B_Y is weakly sequentially compact. Thus there exists $\{x_{n_k}\}$ such that $x_{n_k} \rightarrow x \in B_Y$ weakly in Y . Letting $x^* \in X^*, i^* x^* \in Y'$ and so

$$x^*(x_{n_k}) = i^* x^*(x_{n_k}) \rightarrow i^* x^*(x) = x^*(x)$$

and so in fact, $x_{n_k} \rightarrow x$ weakly in X . ■

The following is the form of the Eberlein Smulian theorem which is often used.

Corollary 21.5.15 *Let $\{x_n\}$ be any bounded sequence in a reflexive Banach space X . Then there exists $x \in X$ and a subsequence, $\{x_{n_k}\}$ such that for all $x^* \in X'$, it follows that $\lim_{k \rightarrow \infty} x^*(x_{n_k}) = x^*(x)$.*

Proof: If a subsequence, x_{n_k} has $\|x_{n_k}\| \rightarrow 0$, then the conclusion follows. Simply let $x = 0$. Suppose then that $\|x_n\|$ is bounded away from 0. That is, $\|x_n\| \in [\delta, C]$. Take a subsequence such that $\|x_{n_k}\| \rightarrow a$. Then consider $x_{n_k}/\|x_{n_k}\|$. By the Eberlein Smulian theorem, this subsequence has a further subsequence, $x_{n_{k_j}}/\|x_{n_{k_j}}\|$ which converges weakly to $x \in B$ where B is the closed unit ball. It follows from routine considerations that $x_{n_{k_j}} \rightarrow ax$ weakly. ■

21.6 Differential Equations

It is a good idea to do Problems 22-24 at this time. Consider $y' = f(t, y, \lambda)$, $y(t_0) = y_0$ for t near t_0 where $\lambda \in V \subseteq \Lambda$ with V an open subset of Λ some Banach space. Assume $f : (t_0 - \delta, t_0 + \delta) \times U \times V \rightarrow Z$ is $C^1((t_0 - \delta, t_0 + \delta) \times U \times V)$ where U is an open subset of Z a Banach space and $u_0 \in U$. Let $\alpha \in (-\delta, \delta)$ and let

$$\alpha s \equiv t - t_0, \phi(s) \equiv y(t) - y_0$$

Thus $\phi(0) = 0$. Also $\phi \in C^1([-1, 1]; Z)$ so $\phi \in \mathcal{D}^1 \equiv \{y \in C^1([-1, 1], Z) : y(0) = 0\}$. From the above problems, \mathcal{D}^1 is a Banach space. It is also the case that

$$\phi'(s) = y'(t)\alpha = \alpha f(\alpha s, y_0 + \phi(s), \lambda)$$

Let L be as in Problem 22, $L\phi = \phi'$. Then the problem reduces to

$$L\phi(s) - \alpha f(\alpha s, y_0 + \phi(s), \lambda) = 0, s \in [-1, 1]$$

Let \mathcal{U}_Z be the open of Problem 24, all $u \in \mathcal{D}^1$ (so $u(0) = 0$) such that $u(t) \in \mathcal{U}_Z$ an open set in Z containing 0, this for each $t \in [-1, 1]$. Let

$$F : (-\delta, \delta) \times U \times \mathcal{U}_Z \times V \rightarrow C([-1, 1]; Z)$$

be defined by

$$F(\alpha, \tilde{y}_0, \psi, \mu)(s) \equiv L\psi(s) - \alpha f(\alpha s, \tilde{y}_0 + \psi(s), \mu)$$

Are the various partial derivatives continuous?

$$\begin{aligned} & F(\alpha + \beta, \tilde{y}_0, \psi, \mu)(s) - F(\alpha, \tilde{y}_0, \psi, \mu)(s) \\ &= \alpha f(\alpha s, \tilde{y}_0 + \psi(s), \mu) - (\alpha + \beta) f((\alpha + \beta)s, \tilde{y}_0 + \psi(s), \mu) \\ &= -\beta f(\alpha s, \tilde{y}_0 + \psi(s), \mu) + (\alpha + \beta) \begin{pmatrix} f(\alpha s, \tilde{y}_0 + \psi(s), \mu) \\ -f((\alpha + \beta)s, \tilde{y}_0 + \psi(s), \mu) \end{pmatrix} \\ &= -\beta f(\alpha s, \tilde{y}_0 + \psi(s), \mu) - (\alpha + \beta) (D_1 f(\alpha s, \tilde{y}_0 + \psi(s), \mu) \beta s + o(\beta s)) \\ &= -\beta f(\alpha s, \tilde{y}_0 + \psi(s), \mu) - \alpha (D_1 f(\alpha s, \tilde{y}_0 + \psi(s), \mu) \beta s + o(\beta s)) \end{aligned}$$

Thus $\alpha \rightarrow D_1 F(\alpha, \tilde{y}_0, \psi, \mu)$ is continuous as a map from $(-\delta, \delta)$ to $\mathcal{L}(\mathbb{R}, C([-1, 1]; Z))$. Similarly, $\tilde{y}_0 \rightarrow D_2 F(\alpha, \tilde{y}_0, \psi, \mu)$ is continuous as a map from U to $\mathcal{L}(Z, C([-1, 1]; Z))$. and $\mu \rightarrow D_4 F(\alpha, \tilde{y}_0, \psi, \mu)$ is continuous as a map from V to $\mathcal{L}(\Lambda, C([-1, 1]; Z))$. What remains is to consider $D_3 F$. Note that $v_n \rightarrow v$ in \mathcal{D}^1 implies $v_n(s) \rightarrow v(s)$ in Z for each s .

$$F(\alpha, \tilde{y}_0, \psi + \eta, \mu)(s) - F(\alpha, \tilde{y}_0, \psi, \mu)(s) + L\eta$$

$$\begin{aligned}
&= F(\alpha, \tilde{y}_0, \psi(s) + \eta(s), \mu) - F(\alpha, \tilde{y}_0, \psi(s), \mu) \\
&= -\alpha D_3 f(\alpha, \tilde{y}_0, \psi(s), \mu) \eta(s) + o(\eta(s)) + L\eta
\end{aligned}$$

Now because of the definition of \mathcal{D}^1 , if $g(\eta)(s) \equiv o(\eta(s))$, then $g(\eta) = o(\eta)$. Indeed, $\frac{\|g(\eta)(s)\|}{\|\eta\|} \leq \frac{\|g(\eta)(s)\|}{\|\eta(s)\|} < \varepsilon$ if $\|\eta\|_{\mathcal{D}^1}$ is small enough. Hence $\|g(\eta)\| \leq \varepsilon \|\eta\|$ if $\|\eta\|$ is small enough. Thus $D_3 F(\alpha, \tilde{y}_0, \psi, \mu)(s) = -\alpha D_3 f(\alpha, \tilde{y}_0, \psi, \mu) + L$, so $\psi \rightarrow D_3 F(\alpha, \tilde{y}_0, \psi, \mu)$ is continuous into $\mathcal{L}(\mathcal{D}^1, C([-1, 1]; Z))$. Since L is continuous, one to one and onto, the open mapping theorem says that $L^{-1} : C([-1, 1]; Z) \rightarrow \mathcal{D}^1$ is continuous, Problem 23. It follows that for given $y_0 \in U$, $D_3 F(0, y_0, 0, \lambda) = L$ which is invertible. By the implicit function theorem, there exists a unique $\phi(\alpha, \tilde{y}_0, \mu)$ where ϕ is a C^1 function of α, \tilde{y}_0, μ for α close enough to 0, \tilde{y}_0 close enough to y_0 and μ close enough to λ , say $(\alpha, \tilde{y}_0, \mu) \in (-\sigma, \sigma) \times B(y_0, r) \times B(\lambda, \delta)$ for which

$$F(\alpha, \tilde{y}_0, \phi, \mu)(s) \equiv L\phi(s) - \alpha f(\alpha s, \tilde{y}_0 + \phi(s, \tilde{y}_0, \mu), \mu) = 0$$

Pick small positive α . Thus, this α will be fixed. Now go backwards in how this started. Let $t = \alpha s + t_0$ so $t - t_0 \in [-\alpha, \alpha]$. Let $y(t, \tilde{y}_0, \mu) \equiv \phi(s, \tilde{y}_0, \mu) + \tilde{y}_0$ so $y'(t) \alpha = \phi'(s) = L\phi = \alpha f(t, y(t, \tilde{y}_0, \mu), \mu)$ so $y(t_0, \tilde{y}_0, \mu) = y_0, y'(t, y_0, \mu) = f(t, y(t, \tilde{y}_0, \mu), \mu)$.

If f had been C^k instead of just C^1 , the same conclusion would follow except now you would have $y(t, \tilde{y}_0, \mu)$ a C^k function of \tilde{y}_0 and μ for (\tilde{y}_0, μ) close to a particular (y_0, λ) .

This use of the implicit function theorem to give existence, uniqueness, and differentiable dependence on initial data and a given parameter is extremely significant because it justifies often used procedures for writing a solution to a differential equation in terms of a Taylor series in powers of a parameter.

This proves the following theorem which is an existence and uniqueness theorem for the initial value problem for ordinary differential equations that also gives a description of dependence on the initial data and an arbitrary parameter.

Theorem 21.6.1 *Let $f : (t_0 - \delta, t_0 + \delta) \times U \times V \rightarrow Z$ where U is an open set in Z, V an open set in Λ , some Banach space. Suppose f is*

$$C^k((t_0 - \delta, t_0 + \delta) \times U \times V, Z).$$

Then if $(y_0, \lambda) \in (U \times V)$, there exists $\alpha > 0$ and a unique solution to

$$y'(t) = f(t, \tilde{y}_0, \mu), y(t_0) = \tilde{y}_0$$

for $t - t_0 \in [-\alpha, \alpha]$ whenever (\tilde{y}_0, μ) is close enough to (y_0, λ) . Denoting this solution as $y(t) = y(t, \tilde{y}_0, \mu)$, it follows that $(\tilde{y}_0, \mu) \rightarrow y(t, \tilde{y}_0, \mu)$ is a C^k function.

The case where t is not just real but is allowed to be complex and f is analytic is also available, but I have not discussed analytic functions here. This leads to being able to expand the solution in a power series. See [11] for more on this subject including the analytic case. This case is also in [36].

Example 21.6.2 *Consider $y' = y^2 + \varepsilon, y(0) = 0$ where ε is a small real number. Then there is a solution $t \rightarrow y(t)$ on a small interval containing 0. According to the above theorem,*

$$y(t) = \sum_{k=0}^2 a_k(t) \varepsilon^k + o(\varepsilon^2)$$

Use the initial condition to find that $a_k(0) = 0$ for each k . Then neglecting higher powers of ε than 2,

$$\sum_{k=0}^2 a'_k(t) \varepsilon^k = 2\varepsilon^2 a_0(t) a_2(t) + \varepsilon^2 a_1^2(t) + 2\varepsilon a_0(t) a_1(t) + a_0^2(t) + \varepsilon$$

Matching the powers of ε , $a'_0(t) = a_0^2(t)$, $a_0(0) = 0$ so $a_0(t) = 0$. Similarly $a'_1(t) = 1$ so $a_1(t) = t$ and $a'_2(t) = t^2$, $a_2(0) = 0$, $a_2(t) = \frac{t^3}{3}$. Thus $y(t) = t\varepsilon + \frac{t^3}{3}\varepsilon^2 + o(\varepsilon^2)$. The exact solution for positive ε is $\sqrt{\varepsilon} \tan(t\sqrt{\varepsilon})$ and this equals $t\varepsilon + \frac{1}{3}t^3\varepsilon^2 + o(\varepsilon^2)$. This is an easy example, but the same idea will hold for harder examples for which it might not be possible to find an exact solution.

21.7 Lyapunov Schmidt Procedure

You have $f : X \times \Lambda \rightarrow Y$ where here X, Λ are Banach spaces and f is C^p . Suppose $(0, 0) \in X \times \Lambda$ and $f(0, 0) = 0$. Then if $D_1 f(0, 0)^{-1}$ is in $\mathcal{L}(Y, X)$, the implicit function theorem says that there exists $x(\lambda)$ a C^p function such that locally $f(x(\lambda), \lambda) = 0$. So what if $D_1 f(0, 0)$ fails to be one to one? Sometimes this case is also considered. It may be that $D_1 f(0, 0)$ is one to one on some subspace and other nice things happen. In particular, suppose the following.

Letting $X_2 \equiv \ker D_1 f(0, 0)$ assume

$$X = X_1 \oplus X_2, \dim(X_2) < \infty$$

where X_1 is a closed subspace. Thus $D_1 f(0, 0)$ is one to one on X_1 . We let

$$Y_1 = D_1 f(0, 0)(X_1)$$

and suppose that $Y = Y_1 \oplus Y_2$ where $\dim(Y_2) < \infty$, and Y_1 is also a closed subspace.

$$\begin{aligned} X_1 &\xrightarrow{D_1 f(0, 0)} Y_1 = D_1 f(0, 0)(X_1), Y_1 \text{ closed} < \infty \\ Y &= Y_1 \oplus Y_2, \dim(Y_2) < \infty \end{aligned}$$

By the open mapping theorem, $D_1 f(0, 0)^{-1}$ is also continuous.

Let Q be a continuous projection onto Y_1 which is assumed to exist² such that $QY_2 = 0$ and $(I - Q)$ is a projection onto Y_2 . Then the equation $f(x(\lambda), \lambda) = 0$ can be written as the pair

$$\begin{aligned} Qf(x, \lambda) &= 0 \\ (I - Q)f(x, \lambda) &= 0 \end{aligned}$$

Consider the top. For $x = x_1 + x_2$ where $x_i \in X_i$, this is

$$Qf(x_1 + x_2, \lambda) = 0$$

Then if $g(x_1, x_2, \lambda) = Qf(x_1 + x_2, \lambda)$, one has $g : X_1 \times X_2 \times \Lambda \rightarrow Y_1$

$$D_1 g(x_1, x_2, \lambda)h = D_1 Qf(x_1 + x_2, \lambda)h, h \in X_1.$$

²In Hilbert space, the existence of this projection map is obvious and it is assumed that it exists here.

Thus $D_1g(0,0,0)^{-1}$ is continuous by the open mapping theorem ($D_1f(0,0)$ is one to one on X_1), and by the implicit function theorem, there is a solution to

$$Qf(x_1 + x_2, \lambda) = 0$$

for $x_1 = x_1(x_2, \lambda)$. (Note how it is important that X_1 and Y_1 be Banach spaces.) Then the other equation yields

$$(I - Q)f(x_1(x_2, \lambda) + x_2, \lambda) = 0$$

and so for fixed λ , this is a finite set of equations of a variable in a finite dimensional space.

This depends on being able to write $X = X_1 \oplus X_2$ where X_1 is closed, $X_2 = \ker D_1f(0,0)$, a similar situation for $Y = Y_1 \oplus Y_2$. So when does this happen? Are there conditions on $D_1f(0,0)$ which will cause it to occur?

There are such conditions. For example, $D_1f(0,0)$ could be a Fredholm operator defined in Definition 21.7.1.

Definition 21.7.1 Let $T \in \mathcal{L}(X, Y)$. Then this is a Fredholm operator means

1. $\dim(\ker(T)) < \infty$
2. $\dim(E) < \infty$ where $Y = TX \oplus E$

The following are some easy examples in which all that nonsense about things being finite dimensional and part of a direct sum does not need to be considered.

Example 21.7.2 Say $X = \mathbb{R}^2$ and $\Lambda = \mathbb{R}$. Let $f(x, y, \lambda) = x + xy + y^2 + \lambda$. Then

$$D_1f(0,0,0) = (1, 0)$$

this 1×2 matrix mapping \mathbb{R}^2 to \mathbb{R} . Thus $X_2 = (0, \alpha)^T : \alpha \in \mathbb{R}$ and $X_1 = (\alpha, 0)^T : \alpha \in \mathbb{R}$. In this case, $Y_1 = \mathbb{R}$ and so $Q = I$. Thus the above reduces to the single equation

$$f((\alpha, 0) + (0, \beta), \lambda) = 0$$

and so since $D_1f(0,0,0)$ is one to one, $x_1 = (\alpha, 0) = x_1((0, \beta), \lambda)$. Of course this is completely obvious because if you consider f in the natural way as a function of three variables, then the implicit function theorem immediately gives $x = x(y, \lambda)$ which is essentially the same result. We just write $(\alpha, 0)$ in place of α . The first independent variable is a function of the other two.

Example 21.7.3 Here is another easy example. $f : \mathbb{R}^2 \times \mathbb{R} \rightarrow \mathbb{R}^2$

$$f(x, y, \lambda) = \begin{pmatrix} x + xy + y^2 + \sin(\lambda) \\ x + y^2 - x^2 + \lambda \end{pmatrix}$$

Then

$$D_1f(x, y, \lambda) = \begin{pmatrix} 1+y & x+2y \\ 1-2x & 2y \end{pmatrix}$$

So

$$D_1f((0,0),0) = \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix}$$

Then

$$X_2 = \ker D_1 \mathbf{f}((0,0),0) = \left\{ \begin{pmatrix} 0 \\ \beta \end{pmatrix} : \beta \in \mathbb{R} \right\}$$

and $X_1 = \left\{ \begin{pmatrix} \alpha \\ 0 \end{pmatrix} : \alpha \in \mathbb{R} \right\}$ and clearly $D_1 \mathbf{f}((0,0),0)$ is indeed one to one on X_1 .

$$D_1 \mathbf{f}(\mathbf{0},0)(X_1) = \left\{ \begin{pmatrix} y \\ y \end{pmatrix} : y \in \mathbb{R} \right\} = Y_1$$

In this case, let

$$Q \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \begin{pmatrix} \frac{\alpha+\beta}{2} \\ \frac{\alpha+\beta}{2} \end{pmatrix} = \begin{pmatrix} 1/2 & 1/2 \\ 1/2 & 1/2 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

so $(I-Q) = \begin{pmatrix} 1/2 & -1/2 \\ -1/2 & 1/2 \end{pmatrix}$. Thus the equations are

$$\begin{aligned} Q \mathbf{f}(\mathbf{x}, \lambda) &= \mathbf{0} \\ (I-Q) \mathbf{f}(\mathbf{x}, \lambda) &= \mathbf{0} \end{aligned}$$

This reduces to

$$\begin{aligned} \begin{pmatrix} -\frac{1}{2}x^2 + \frac{1}{2}xy + x + y^2 + \frac{1}{2}\lambda + \frac{1}{2}\sin \lambda \\ -\frac{1}{2}x^2 + \frac{1}{2}xy + x + y^2 + \frac{1}{2}\lambda + \frac{1}{2}\sin \lambda \end{pmatrix} &= \begin{pmatrix} 0 \\ 0 \end{pmatrix} \\ \begin{pmatrix} \frac{1}{2}x^2 + \frac{1}{2}yx - \frac{1}{2}\lambda + \frac{1}{2}\sin \lambda \\ -\frac{1}{2}x^2 - \frac{1}{2}yx + \frac{1}{2}\lambda - \frac{1}{2}\sin \lambda \end{pmatrix} &= \begin{pmatrix} 0 \\ 0 \end{pmatrix} \end{aligned}$$

Note how in both the top and the bottom, there is only one equation and one can solve for x in terms of y, λ near $(0,0,0)$ which is what the above general argument shows. Of course you can see this directly using the implicit function theorem. Then can you solve for $y = y(\lambda)$? This would involve trying to solve for y as a function of λ in the following where $x(y, \lambda)$ comes from the first equations.

$$\frac{1}{2}x^2(y, \lambda) + \frac{1}{2}yx(y, \lambda) - \frac{1}{2}\lambda + \frac{1}{2}\sin \lambda = 0$$

If you can do this, then you would have found (x, y) as a function of λ for small λ .

In this example, in the top equation, at $(0,0,0)$, $x_y = 0$. Also $x_\lambda = -1$ so $x(y, \lambda) \approx -\lambda$ other than higher order terms for small y, λ . Then in the bottom equation, for all variables very small, you would have $\lambda^2 + y(-\lambda) - \lambda + \sin(\lambda) = 0$, $y(\lambda) = -1 + \frac{\sin(\lambda)}{\lambda} + \lambda$ at least approximately. Thus it seems there is a nonzero solution to the equation $\mathbf{f}(x, y, \lambda) = \mathbf{0}$ which is valid for small λ, x, y , this in addition to the zero solution. Note that for small nonzero λ , $-1 + \frac{\sin(\lambda)}{\lambda} + \lambda \neq 0$. It equals approximately $\lambda - \frac{\lambda^2}{3!}$ for small λ from the power series for \sin .

In the next example, the same procedure gives a solution to a problem

$$\mathbf{f}((x, y), \lambda) = \mathbf{0}$$

such that for small λ , (x, y) is a function of λ which is nonzero and

$$\mathbf{f}((0,0), \lambda) = \mathbf{0}$$

Thus for small λ , there are two solutions to the nonlinear system of equations.

Example 21.7.4 *Let*

$$f((x, y), \lambda) = \begin{pmatrix} x + xy + y^2 + x \sin(\lambda) \\ x + y^2 - x^2 + x\lambda \end{pmatrix}$$

In this case $f((0, 0), \lambda) = \mathbf{0}$ even though λ might not be 0. The Lyapunov Schmidt procedure will be used to show that there are nonzero solutions $x(\lambda), y(\lambda)$ such that

$$f((x(\lambda), y(\lambda)), \lambda) = \mathbf{0}$$

At origin,

$$D_1 f((0, 0), 0) = \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix}$$

Thus $X_1 = \text{span}(e_1)$ and $X_2 = \text{span}(e_2)$. Then $Y_1 = \text{span}(e_1 + e_2)$ and $Y_2 = \text{span}(e_1 - e_2)$. Also $D_1 f((0, 0), 0)$ is one to one on X_1 and its range is Y_1 . Then let

$$Q \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \begin{pmatrix} \frac{\alpha+\beta}{2} \\ \frac{\alpha+\beta}{2} \end{pmatrix} = \begin{pmatrix} 1/2 & 1/2 \\ 1/2 & 1/2 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

$$(I - Q) = \begin{pmatrix} 1/2 & -1/2 \\ -1/2 & 1/2 \end{pmatrix}$$

Then $Qf = \mathbf{0}$ is yields the equation

$$x + \frac{1}{2}x\lambda + \frac{1}{2}x\sin\lambda + \frac{1}{2}xy - \frac{1}{2}x^2 + y^2 = 0$$

Also $(I - Q)f = \mathbf{0}$ yields the equation

$$\frac{1}{2}x\sin\lambda - \frac{1}{2}x\lambda + \frac{1}{2}xy + \frac{1}{2}x^2 = 0$$

Now consider x_y and x_λ at $(0, 0)$ from the first equation. Both of these are easily seen to be 0. Now consider x_{yy} . After some computations, this is seen to be $x_{yy} = -2$. Similarly, $x_{y\lambda}(0, 0) = 0, x_{\lambda\lambda}(0, 0) = 0$ also. Thus up to terms of degree 3,

$$x(y, \lambda) = -y^2 = \frac{1}{2}(-2)y^2$$

Place this in the bottom equation.

$$\frac{1}{2}y^2\lambda - \frac{1}{2}y^2\sin\lambda - \frac{1}{2}y^3 + \frac{1}{2}y^4 = 0$$

Now the idea is to find $y = y(\lambda)$, hopefully nonzero. Divide by y^2 and multiply by 2.

$$y^2 - y + \lambda - \sin\lambda = 0$$

Then for small λ this is approximately equal to

$$y^2 - y + \frac{\lambda^3}{6} = 0$$

Then a solution for y for small λ is

$$y = \frac{1 + \sqrt{1 - \frac{2}{3}\lambda^3}}{2}$$

Of course there is another solution as well, when you replace the $+$ with a minus sign. This is the one we want because when $\lambda = 0$ it reduces to $y = 0$. This shows that there exist solutions to the equations $\mathbf{f}((x, y), \lambda) = \mathbf{0}$ which for small λ are approximately

$$(x(\lambda), y(\lambda)) = \left(-y^2, \frac{1 - \sqrt{1 - \frac{2}{3}\lambda^3}}{2} \right)$$

In terms of λ very small,

$$(x(\lambda), y(\lambda)) = \left(\frac{1}{6}\lambda^3 + \frac{1}{6}\sqrt{3}\sqrt{3 - 2\lambda^3} - \frac{1}{2}, \frac{1 - \sqrt{1 - \frac{2}{3}\lambda^3}}{2} \right)$$

Using a power series in λ to approximate these functions, this reduces to

$$(x(\lambda), y(\lambda)) = \left(-\frac{1}{36}\lambda^6, \frac{1}{6}\lambda^3 + \frac{1}{36}\lambda^6 + \frac{1}{108}\lambda^9 \right)$$

where higher order terms are neglected. Thus there exist other solutions than the zero solution even though λ may be nonzero. Note that in this example, $\mathbf{f}((0, 0), \lambda) = \mathbf{0}$.

Note that all of this works as well if the function f is defined on an open subset of $X \times \Lambda$ because it is really just an application of the implicit function theorem.

21.8 The Holder Spaces

This is such an important example that I am including it. It is an example of a Banach space which is not separable.

Definition 21.8.1 Let $p > 1$. Then $f \in C^{1/p}([0, 1])$ means that $f \in C([0, 1])$ and also

$$\rho_p(f) \equiv \sup \left\{ \frac{|f(x) - f(y)|}{|x - y|^{1/p}} : x, y \in X, x \neq y \right\} < \infty$$

Then the norm is defined as $\|f\|_{C([0, 1])} + \rho_p(f) \equiv \|f\|_{1/p}$.

It is an exercise to verify that $C^{1/p}([0, 1])$ is a Banach space.

Let $p > 1$. Then $C^{1/p}([0, 1])$ is not separable. Define uncountably many functions, one for each ε where ε is a sequence of -1 and 1 . Thus $\varepsilon_k \in \{-1, 1\}$. Thus $\varepsilon \neq \varepsilon'$ if the two sequences differ in at least one slot, one giving 1 and the other equaling -1 . There are uncountably many of these sequences, equal to the number of subsets of \mathbb{N} . Now define $f_\varepsilon(t) \equiv \sum_{k=1}^{\infty} \varepsilon_k 2^{-k/p} \sin(2^k \pi t)$. Then this is $1/p$ Holder. Let $s < t$.

$$|f_\varepsilon(t) - f_\varepsilon(s)| \leq \sum_{k \leq |\log_2(t-s)|} \left| 2^{-k/p} \sin(2^k \pi t) - 2^{-k/p} \sin(2^k \pi s) \right|$$

$$+ \sum_{k > |\log_2(t-s)|} \left| 2^{-k/p} \sin(2^k \pi t) - 2^{-k/p} \sin(2^k \pi s) \right|$$

If $t = 1$ and $s = 0$, there is really nothing to show because then the difference equals 0. There is also nothing to show if $t = s$. From now on, $0 < t - s < 1$. Let k_0 be the largest integer which is less than or equal to $|\log_2(t-s)| = -\log_2(t-s)$. Note that $-\log(t-s) > 0$ because $0 < t-s < 1$. Then

$$\begin{aligned} |f_\varepsilon(t) - f_\varepsilon(s)| &\leq \sum_{k \leq k_0} \left| 2^{-k/p} \sin(2^k \pi t) - 2^{-k/p} \sin(2^k \pi s) \right| \\ &\quad + \sum_{k > k_0} \left| 2^{-k/p} \sin(2^k \pi t) - 2^{-k/p} \sin(2^k \pi s) \right| \\ &\leq \sum_{k \leq k_0} 2^{-k/p} 2^k \pi |t-s| + \sum_{k > k_0} 2^{-k/p} 2 \end{aligned}$$

Now $k_0 \leq -\log_2(t-s) < k_0 + 1$ and so $-k_0 \geq \log_2(t-s) \geq -(k_0 + 1)$. Hence $2^{-k_0} \geq |t-s| \geq 2^{-k_0} 2^{-1}$ and so $2^{-k_0/p} \geq |t-s|^{1/p} \geq 2^{-k_0/p} 2^{-1/p}$. Using this in the sums,

$$\begin{aligned} |f_\varepsilon(t) - f_\varepsilon(s)| &\leq |t-s| C_p + \sum_{k > k_0} 2^{-k/p} 2^{k_0/p} 2^{-k_0/p} 2 \\ &\leq |t-s| C_p + \sum_{k > k_0} 2^{-k/p} 2^{k_0/p} \left(2^{1/p} |t-s|^{1/p} \right) 2 \\ &\leq |t-s| C_p + \sum_{k > k_0} 2^{-(k-k_0)/p} \left(2^{1/p} |t-s|^{1/p} \right) 2 \\ &\leq C_p |t-s| + \left(2^{1+1/p} \right) \sum_{k=1}^{\infty} 2^{-k/p} |t-s|^{1/p} \\ &= C_p |t-s| + D_p |t-s|^{1/p} \leq C_p |t-s|^{1/p} + D_p |t-s|^{1/p} \end{aligned}$$

Thus f_ε is indeed $1/p$ Holder continuous.

Now consider $\varepsilon \neq \varepsilon'$. Suppose the first discrepancy in the two sequences occurs with ε_j . Thus one is 1 and the other is -1 . Let $t = \frac{i+1}{2^{j+1}}, s = \frac{i}{2^{j+1}}$

$$\begin{aligned} |f_\varepsilon(t) - f_\varepsilon(s) - (f_{\varepsilon'}(t) - f_{\varepsilon'}(s))| &= \\ \left| \sum_{k=j}^{\infty} \varepsilon_k 2^{-k/p} \sin(2^k \pi t) - \sum_{k=j}^{\infty} \varepsilon_k 2^{-k/p} \sin(2^k \pi s) \right. \\ \left. - \left(\sum_{k=j}^{\infty} \varepsilon'_k 2^{-k/p} \sin(2^k \pi t) - \sum_{k=j}^{\infty} \varepsilon'_k 2^{-k/p} \sin(2^k \pi s) \right) \right| \end{aligned}$$

Now consider what happens for $k > j$. Then $\sin\left(2^k \pi \frac{i}{2^{j+1}}\right) = \sin(m\pi) = 0$ for some integer m . Thus the whole mess reduces to

$$\begin{aligned} &\left| (\varepsilon_j - \varepsilon'_j) 2^{-j/p} \sin\left(\frac{2^j \pi (i+1)}{2^{j+1}}\right) - (\varepsilon_j - \varepsilon'_j) 2^{-j/p} \sin\left(\frac{2^j \pi i}{2^{j+1}}\right) \right| \\ &= \left| (\varepsilon_j - \varepsilon'_j) 2^{-j/p} \sin\left(\frac{\pi(i+1)}{2}\right) - (\varepsilon_j - \varepsilon'_j) 2^{-j/p} \sin\left(\frac{\pi i}{2}\right) \right| \\ &= 2 \left(2^{-j/p} \right) \end{aligned}$$

In particular, $|t - s| = \frac{1}{2^{j+1}}$ so $2^{1/p} |t - s|^{1/p} = 2^{-j/p}$

$$|f_{\epsilon}(t) - f_{\epsilon}(s) - (f_{\epsilon'}(t) - f_{\epsilon'}(s))| = 2 \left(2^{1/p} \right) |t - s|^{1/p}$$

which shows that

$$\sup_{0 \leq s < t \leq 1} \frac{|f_{\epsilon}(t) - f_{\epsilon'}(t) - (f_{\epsilon}(s) - f_{\epsilon'}(s))|}{|t - s|^{1/p}} \geq 2^{1/p} (2)$$

Thus there exists a set of uncountably many functions in $C^{1/p}([0, T])$ and for any two of them f, g , you get

$$\|f - g\|_{C^{1/p}([0,1])} > 2$$

so $C^{1/p}([0, 1])$ is not separable.

21.9 Exercises

1. Is \mathbb{N} a G_{δ} set? What about \mathbb{Q} ? What about a countable dense subset of a complete metric space?
2. \uparrow Let $f : \mathbb{R} \rightarrow \mathbb{C}$ be a function. Define the oscillation of a function in $B(x, r)$ by $\omega_r f(x) = \sup\{|f(z) - f(y)| : y, z \in B(x, r)\}$. Define the oscillation of the function at the point, x by $\omega f(x) = \lim_{r \rightarrow 0} \omega_r f(x)$. Show f is continuous at x if and only if $\omega f(x) = 0$. Then show the set of points where f is continuous is a G_{δ} set (try $U_n = \{x : \omega f(x) < \frac{1}{n}\}$). Does there exist a function continuous at only the rational numbers? Does there exist a function continuous at every irrational and discontinuous elsewhere? **Hint:** Suppose D is any countable set, $D = \{d_i\}_{i=1}^{\infty}$, and define the function, $f_n(x)$ to equal zero for every $x \notin \{d_1, \dots, d_n\}$ and 2^{-n} for x in this finite set. Then consider $g(x) \equiv \sum_{n=1}^{\infty} f_n(x)$. Show that this series converges uniformly.
3. Let $f \in C([0, 1])$ and suppose $f'(x)$ exists. Show there exists a constant, K , such that $|f(x) - f(y)| \leq K|x - y|$ for all $y \in [0, 1]$. Let $U_n = \{f \in C([0, 1]) \text{ such that for each } x \in [0, 1] \text{ there exists } y \in [0, 1] \text{ such that } |f(x) - f(y)| > n|x - y|\}$. Show that U_n is open and dense in $C([0, 1])$ where for $f \in C([0, 1])$,

$$\|f\| \equiv \sup\{|f(x)| : x \in [0, 1]\}.$$

Show that $\cap_n U_n$ is a dense G_{δ} set of nowhere differentiable continuous functions. Thus every continuous function is uniformly close to one which is nowhere differentiable.

4. \uparrow Suppose $f(x) = \sum_{k=1}^{\infty} u_k(x)$ where the convergence is uniform and each u_k is a polynomial. Is it reasonable to conclude that $f'(x) = \sum_{k=1}^{\infty} u'_k(x)$? The answer is no. Use Problem 3 and the Weierstrass approximation theorem to show this.
5. Let X be a normed linear space. $A \subseteq X$ is “weakly bounded” if for each $x^* \in X'$, $\sup\{|x^*(x)| : x \in A\} < \infty$, while A is bounded if $\sup\{\|x\| : x \in A\} < \infty$. Show A is weakly bounded if and only if it is bounded.

6. Let f be a 2π periodic locally integrable function on \mathbb{R} . The Fourier series for f is given by

$$\sum_{k=-\infty}^{\infty} a_k e^{ikx} \equiv \lim_{n \rightarrow \infty} \sum_{k=-n}^n a_k e^{ikx} \equiv \lim_{n \rightarrow \infty} S_n f(x)$$

where

$$a_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{-ikx} f(x) dx.$$

Show

$$S_n f(x) = \int_{-\pi}^{\pi} D_n(x-y) f(y) dy$$

where

$$D_n(t) = \frac{\sin((n + \frac{1}{2})t)}{2\pi \sin(\frac{t}{2})}.$$

Verify that $\int_{-\pi}^{\pi} D_n(t) dt = 1$. Also show that if $g \in L^1(\mathbb{R})$, then

$$\lim_{a \rightarrow \infty} \int_{\mathbb{R}} g(x) \sin(ax) dx = 0.$$

This last is called the Riemann Lebesgue lemma. **Hint:** For the last part, assume first that $g \in C_c^\infty(\mathbb{R})$ and integrate by parts. Then exploit density of the set of functions in $L^1(\mathbb{R})$.

7. \uparrow It turns out that the Fourier series sometimes converges to the function pointwise. Suppose f is 2π periodic and Holder continuous. That is $|f(x) - f(y)| \leq K|x - y|^\theta$ where $\theta \in (0, 1]$. Show that if f is like this, then the Fourier series converges to f at every point. Next modify your argument to show that if at every point, x , $|f(x+) - f(y)| \leq K|x - y|^\theta$ for y close enough to x and larger than x and

$$|f(x-) - f(y)| \leq K|x - y|^\theta$$

for every y close enough to x and smaller than x , then $S_n f(x) \rightarrow \frac{f(x+) + f(x-)}{2}$, the midpoint of the jump of the function. **Hint:** Use Problem 6.

8. \uparrow Let $Y = \{f \text{ such that } f \text{ is continuous, defined on } \mathbb{R}, \text{ and } 2\pi \text{ periodic}\}$. Define $\|f\|_Y = \sup\{|f(x)| : x \in [-\pi, \pi]\}$. Show that $(Y, \|\cdot\|_Y)$ is a Banach space. Let $x \in \mathbb{R}$ and define $L_n(f) = S_n f(x)$. Show $L_n \in Y'$ but $\lim_{n \rightarrow \infty} \|L_n\| = \infty$. Show that for each $x \in \mathbb{R}$, there exists a dense G_δ subset of Y such that for f in this set, $|S_n f(x)|$ is unbounded. Finally, show there is a dense G_δ subset of Y having the property that $|S_n f(x)|$ is unbounded on the rational numbers. **Hint:** To do the first part, let $f(y)$ approximate $\text{sgn}(D_n(x - y))$. Here $\text{sgn } r = 1$ if $r > 0$, -1 if $r < 0$ and 0 if $r = 0$. This rules out one possibility of the uniform boundedness principle. After this, show the countable intersection of dense G_δ sets must also be a dense G_δ set.

9. Let $\alpha \in (0, 1]$. Define, for X a compact subset of \mathbb{R}^p ,

$$C^\alpha(X; \mathbb{R}^n) \equiv \{f \in C(X; \mathbb{R}^n) : \rho_\alpha(f) + \|f\| \equiv \|f\|_\alpha < \infty\}$$

where

$$\|f\| \equiv \sup\{|f(x)| : x \in X\}$$

and

$$\rho_\alpha(f) \equiv \sup\left\{\frac{|f(x) - f(y)|}{|x - y|^\alpha} : x, y \in X, x \neq y\right\}.$$

Show that $(C^\alpha(X; \mathbb{R}^n), \|\cdot\|_\alpha)$ is a complete normed linear space. This is called a Holder space. What would this space consist of if $\alpha > 1$?

10. †Let X be the Holder functions which are periodic of period 2π . Define $L_n f(x) = S_n f(x)$ where $L_n : X \rightarrow Y$ for Y given in Problem 8. Show $\|L_n\|$ is bounded independent of n . Conclude that $L_n f \rightarrow f$ in Y for all $f \in X$. In other words, for the Holder continuous and 2π periodic functions, the Fourier series converges to the function uniformly. **Hint:** $L_n f(x)$ is given by

$$L_n f(x) = \int_{-\pi}^{\pi} D_n(y) f(x-y) dy$$

where $f(x-y) = f(x) + g(x, y)$ where $|g(x, y)| \leq C|y|^\alpha$. Use the fact the Dirichlet kernel integrates to one to write

$$\begin{aligned} \left| \int_{-\pi}^{\pi} D_n(y) f(x-y) dy \right| &\leq \overbrace{\left| \int_{-\pi}^{\pi} D_n(y) f(x) dy \right|}^{=|f(x)|} \\ &+ C \left| \int_{-\pi}^{\pi} \sin\left(\left(n + \frac{1}{2}\right)y\right) (g(x, y) / \sin(y/2)) dy \right| \end{aligned}$$

Show the functions, $y \rightarrow g(x, y) / \sin(y/2)$ are bounded in L^1 independent of x and get a uniform bound on $\|L_n\|$. Now use a similar argument to show $\{L_n f\}$ is equicontinuous in addition to being uniformly bounded. In doing this you might proceed as follows. Show

$$\begin{aligned} |L_n f(x) - L_n f(x')| &\leq \left| \int_{-\pi}^{\pi} D_n(y) (f(x-y) - f(x'-y)) dy \right| \\ &\leq \|f\|_\alpha |x - x'|^\alpha \\ &+ \left| \int_{-\pi}^{\pi} \sin\left(\left(n + \frac{1}{2}\right)y\right) \left(\frac{f(x-y) - f(x) - (f(x'-y) - f(x'))}{\sin(y/2)} \right) dy \right| \end{aligned}$$

Then split this last integral into two cases, one for $|y| < \eta$ and one where $|y| \geq \eta$. If $L_n f$ fails to converge to f uniformly, then there exists $\varepsilon > 0$ and a subsequence, n_k such that $\|L_{n_k} f - f\|_\infty \geq \varepsilon$ where this is the norm in Y or equivalently the sup norm on $[-\pi, \pi]$. By the Arzela Ascoli theorem, there is a further subsequence, $L_{n_{k_l}} f$ which converges uniformly on $[-\pi, \pi]$. But by Problem 7 $L_n f(x) \rightarrow f(x)$.

11. Let X be a normed linear space and let M be a convex open set containing 0. Define

$$\rho(x) = \inf\{t > 0 : \frac{x}{t} \in M\}.$$

Show ρ is a gauge function defined on X . This particular example is called a Minkowski functional. It is of fundamental importance in the study of locally convex topological vector spaces. A set, M , is convex if $\lambda x + (1 - \lambda)y \in M$ whenever $\lambda \in [0, 1]$ and $x, y \in M$.

12. \uparrow The Hahn Banach theorem can be used to establish separation theorems. Let M be an open convex set containing 0. Let $x \notin M$. Show there exists $x^* \in X'$ such that $\operatorname{Re} x^*(x) \geq 1 > \operatorname{Re} x^*(y)$ for all $y \in M$. **Hint:** If $y \in M, \rho(y) < 1$. Show this. If $x \notin M, \rho(x) \geq 1$. Try $f(\alpha x) = \alpha \rho(x)$ for $\alpha \in \mathbb{R}$. Then extend f to the whole space using the Hahn Banach theorem and call the result F , show F is continuous, then fix it so F is the real part of $x^* \in X'$.
13. A Banach space is said to be strictly convex if whenever $\|x\| = \|y\|$ and $x \neq y$, then

$$\left\| \frac{x+y}{2} \right\| < \|x\|.$$

$F : X \rightarrow X'$ is said to be a duality map if it satisfies the following: a.) $\|F(x)\| = \|x\|$. b.) $F(x)(x) = \|x\|^2$. Show that if X' is strictly convex, then such a duality map exists. The duality map is an attempt to duplicate some of the features of the Riesz map in Hilbert space. This Riesz map is the map which takes a Hilbert space to its dual defined as follows.

$$R(x)(y) = (y, x)$$

The Riesz representation theorem for Hilbert space says this map is onto. **Hint:** For an arbitrary Banach space, let

$$F(x) \equiv \left\{ x^* : \|x^*\| \leq \|x\| \text{ and } x^*(x) = \|x\|^2 \right\}$$

Show $F(x) \neq \emptyset$ by using the Hahn Banach theorem on $f(\alpha x) = \alpha \|x\|^2$. Next show $F(x)$ is closed and convex. Finally show that you can replace the inequality in the definition of $F(x)$ with an equal sign. Now use strict convexity to show there is only one element in $F(x)$.

14. Prove the following theorem which is an improved version of the open mapping theorem, [14]. Let X and Y be Banach spaces and let $A \in \mathcal{L}(X, Y)$. Then the following are equivalent.

$$AX = Y,$$

A is an open map.

Note this gives the equivalence between A being onto and A being an open map. The open mapping theorem says that if A is onto then it is open.

15. Suppose $D \subseteq X$ and D is dense in X . Suppose $L : D \rightarrow Y$ is linear and $\|Lx\| \leq K\|x\|$ for all $x \in D$. Show there is a unique extension of L, \tilde{L} , defined on all of X with $\|\tilde{L}x\| \leq K\|x\|$ and \tilde{L} is linear. You do not get uniqueness when you use the Hahn Banach theorem. Therefore, in the situation of this problem, it is better to use this result.
16. \uparrow A Banach space is uniformly convex if whenever $\|x_n\|, \|y_n\| \leq 1$ and $\|x_n + y_n\| \rightarrow 2$, it follows that $\|x_n - y_n\| \rightarrow 0$. Show uniform convexity implies strict convexity (See Problem 13). **Hint:** Suppose it is not strictly convex. Then there exist $\|x\|$ and $\|y\|$ both equal to 1 and $\left\| \frac{x_n + y_n}{2} \right\| = 1$ consider $x_n \equiv x$ and $y_n \equiv y$, and use the conditions for uniform convexity to get a contradiction. It can be shown that L^p is uniformly convex whenever $\infty > p > 1$. See Hewitt and Stromberg [26] or Ray [47]. See Theorem 21.3.10.

17. Show that a closed subspace of a reflexive Banach space is reflexive. This is done in the chapter. However, try to do it yourself.
18. x_n converges weakly to x if for every $x^* \in X'$, $x^*(x_n) \rightarrow x^*(x)$. $x_n \rightharpoonup x$ denotes weak convergence. Show that if $\|x_n - x\| \rightarrow 0$, then $x_n \rightharpoonup x$.
19. \uparrow Show that if X is uniformly convex, then if $x_n \rightharpoonup x$ and $\|x_n\| \rightarrow \|x\|$, it follows $\|x_n - x\| \rightarrow 0$. **Hint:** Use Lemma 21.2.9 to obtain $f \in X'$ with $\|f\| = 1$ and $f(x) = \|x\|$. See Problem 16 for the definition of uniform convexity. Now by the weak convergence, you can argue that if $x \neq 0$, $f(x_n/\|x_n\|) \rightarrow f(x/\|x\|)$. You also might try to show this in the special case where $\|x_n\| = \|x\| = 1$.
20. Suppose $L \in \mathcal{L}(X, Y)$ and $M \in \mathcal{L}(Y, Z)$. Show $ML \in \mathcal{L}(X, Z)$ and that $(ML)^* = L^*M^*$.
21. This problem gives a simple condition for the subgradient of a convex function to be onto. Let X be a reflexive Banach space and suppose $\phi : X \rightarrow (-\infty, \infty]$ is convex, proper, lower semicontinuous. This means
- (a) Convex: $\phi(tx + (1-t)y) \leq t\phi(x) + (1-t)\phi(y)$ for all $t \in [0, 1]$.
 - (b) Lower semicontinuous: If $x_n \rightarrow x$, then $\phi(x) \leq \liminf_{n \rightarrow \infty} \phi(x_n)$.
 - (c) The subgradient of ϕ at x , denoted as $\partial\phi(x)$ is defined as follows: $y^* \in \partial\phi(x)$ means $y^*(z - x) \leq \phi(z) - \phi(x)$ for any z .

Suppose then that for all $y^* \in X'$,

$$\lim_{\|x\| \rightarrow \infty} \phi(x) - \langle y^*, x \rangle = \infty$$

this last condition being called “coercive”. Show that under these conditions, you can conclude that $\partial\phi$ is not just nonempty for some x but that in fact every $y^* \in X'$ is contained in some $\partial\phi(x)$. Thus $\partial\phi$ is actually onto. **Hint:** Consider the function $x \rightarrow \phi(x) - \langle y^*, x \rangle$. Argue that it is lower semicontinuous. Let

$$\lambda \equiv \inf \{ \phi(x) - \langle y^*, x \rangle : x \in X \}$$

Let $\{x_n\}$ be a minimizing sequence. Argue that from the coercivity condition, $\|x_n\|$ must be bounded. Now use the Eberlein Smulian theorem, to verify that there is a weakly convergent subsequence $x_n \rightarrow x$ weakly. In finite dimensions, you just use the Heine Borel theorem. You know the epigraph of ϕ intersected with $X \times \mathbb{R}$ is a convex and closed subset of $X \times \mathbb{R}$. Explain why this is so. This will require a separation theorem in infinite dimensional space like Problem 12 above. In finite dimensional space, there isn't much to show here. Next explain why ϕ must be weakly lower semicontinuous. If you can't do this part, just use the theorem that a function which is convex and lower semicontinuous is also weakly lower semicontinuous or specialize to finite dimensions and use advanced calculus. That is, if $x_n \rightarrow x$ weakly, then $\phi(x) \leq \liminf_{n \rightarrow \infty} \phi(x_n)$. Conclude that $\lambda > -\infty$ and equals $\phi(x) - \langle y^*, x \rangle$ which is no larger than $\phi(z) - \langle y^*, z \rangle$. Now conclude that $y^* \in \partial\phi(x)$.

22. Let Z be a Banach space. Let $\mathcal{D}^1 \equiv \{y \in C^1([-1, 1], Z) : y(0) = 0\}$. Let $\|y\|_{\mathcal{D}^1} \equiv \max(\|y\|_\infty, \|y'\|_\infty)$. Show \mathcal{D}^1 is a Banach space.

23. †Let $L : \mathcal{D}^1 \rightarrow C([-1, 1], Z)$, $Ly \equiv y'$. Show L is continuous, defined on \mathcal{D}^1 and one to one onto $L(\mathcal{D}^1)$. Show $L(\mathcal{D}^1) = C([-1, 1], Z)$. Thus L^{-1} is continuous by the open mapping theorem. **Hint:** Adapt the Riemann integral to an integral which has values in a Banach space including the fundamental theorem of calculus. Then if $u \in C([-1, 1], Z)$, consider $\int_0^t u(s) ds \equiv w(t)$ and argue $Lw = u$.
24. †Let U_Z denote an open set in Z . Let $f : U_Z \rightarrow Z$ be C^1 . Then define for $u \in \mathcal{D}^1$, $f(u)(t) \equiv f(u(t))$. Show that $f(\mathcal{D}^1) \subseteq C([-1, 1], Z)$. If \mathcal{U}_Z consists of $u \in \mathcal{D}^1$ such that $u(t) \in U_Z$ for each $t \in [-1, 1]$, show that \mathcal{U}_Z is an open subset of \mathcal{D}^1 . If $f : U \rightarrow Z$ is C^1 , show that $f : \mathcal{U}_Z \rightarrow \mathcal{D}^1$ is also C^1 and that

$$Df(u)(v)(t) = Df(u(t))(v(t)).$$

Chapter 22

Hilbert Spaces

In this chapter, Hilbert spaces, which have been alluded to earlier are given a complete discussion. These spaces, as noted earlier are just complete inner product spaces.

22.1 Basic Theory

Definition 22.1.1 Let X be a vector space. An inner product is a mapping from $X \times X$ to \mathbb{C} if X is complex and from $X \times X$ to \mathbb{R} if X is real, denoted by (x, y) which satisfies the following.

$$(x, x) \geq 0, (x, x) = 0 \text{ if and only if } x = 0, \quad (22.1)$$

$$(x, y) = \overline{(y, x)}. \quad (22.2)$$

For $a, b \in \mathbb{C}$ and $x, y, z \in X$,

$$(ax + by, z) = a(x, z) + b(y, z). \quad (22.3)$$

Note that 22.2 and 22.3 imply $(x, ay + bz) = \bar{a}(x, y) + \bar{b}(x, z)$. Such a vector space is called an inner product space.

The Cauchy Schwarz inequality is fundamental for the study of inner product spaces.

Theorem 22.1.2 (Cauchy Schwarz) In any inner product space

$$|(x, y)| \leq \|x\| \|y\|.$$

Equality holds if and only if x is a multiple of y . The inequality holds under the weaker assumption that $(x, x) \geq 0$ without the stipulation that this happens only if $x = 0$.

Proof: Let $\omega \in \mathbb{C}$, $|\omega| = 1$, and $\bar{\omega}(x, y) = |(x, y)| = \operatorname{Re}(x, y\omega)$. Let $F(t) = (x + ty\omega, x + ty\omega)$, $t \in \mathbb{R}$. If $y = 0$ there is nothing to prove because $(x, 0) = (x, 0 + 0) = (x, 0) + (x, 0)$ and so $(x, 0) = 0$. Thus, it can be assumed $y \neq 0$. Then from the axioms of the inner product,

$$F(t) = \|x\|^2 + 2t \operatorname{Re}(x, \omega y) + t^2 \|y\|^2 \geq 0.$$

This yields $\|x\|^2 + 2t|(x, y)| + t^2 \|y\|^2 \geq 0$. Since this inequality holds for all $t \in \mathbb{R}$, it follows from the quadratic formula that $4|(x, y)|^2 - 4\|x\|^2 \|y\|^2 \leq 0$. In getting this inequality, it was only necessary to assume $(x, x) \geq 0$.

Consider the claim about equality. Note that if $x = \alpha y$, then equality holds directly. Indeed this is the only way this can happen because if equality holds in the Cauchy Schwarz inequality, $F(t) = 0$ for some real t . This happens only if $x = -t\omega y$ for some t real. ■

Proposition 22.1.3 For an inner product space, $\|x\| \equiv (x, x)^{1/2}$ does specify a norm.

Proof: All the axioms are obvious except the triangle inequality. To verify this,

$$\begin{aligned} \|x + y\|^2 &\equiv (x + y, x + y) \equiv \|x\|^2 + \|y\|^2 + 2 \operatorname{Re}(x, y) \\ &\leq \|x\|^2 + \|y\|^2 + 2|(x, y)| \\ &\leq \|x\|^2 + \|y\|^2 + 2\|x\| \|y\| = (\|x\| + \|y\|)^2. \quad \blacksquare \end{aligned}$$

The following lemma is called the parallelogram identity.

Lemma 22.1.4 *In an inner product space,*

$$\|x+y\|^2 + \|x-y\|^2 = 2\|x\|^2 + 2\|y\|^2.$$

The proof, a straightforward application of the inner product axioms, is left to the reader.

Lemma 22.1.5 *For $x \in H$, an inner product space,*

$$\|x\| = \sup_{\|y\| \leq 1} |(x, y)| \quad (22.4)$$

Proof: By the Cauchy Schwarz inequality, if $x \neq 0$,

$$\|x\| \geq \sup_{\|y\| \leq 1} |(x, y)| \geq \left(x, \frac{x}{\|x\|} \right) = \|x\|.$$

It is obvious that 22.4 holds in the case that $x = 0$.

Definition 22.1.6 *A Hilbert space is an inner product space which is complete. Thus a Hilbert space is a Banach space in which the norm comes from an inner product as described above.*

In Hilbert space, one can define a projection map onto closed convex nonempty sets.

Definition 22.1.7 *A set K is convex if whenever $\lambda \in [0, 1]$ and $x, y \in K$, $\lambda x + (1 - \lambda)y \in K$.*

Theorem 22.1.8 *Let K be a closed convex nonempty subset of a Hilbert space H , and let $x \in H$. Then there exists a unique point $Px \in K$ such that $\|Px - x\| \leq \|y - x\|$ for all $y \in K$.*

Proof: Consider uniqueness. Suppose that z_1 and z_2 are two different elements of K such that for $i = 1, 2$,

$$\|z_i - x\| \leq \|y - x\| \quad (22.5)$$

for all $y \in K$. Also, note that since K is convex, $\frac{z_1 + z_2}{2} \in K$. Therefore, by the parallelogram identity,

$$\begin{aligned} \|z_1 - x\|^2 &\leq \left\| \frac{z_1 + z_2}{2} - x \right\|^2 = \left\| \frac{z_1 - x}{2} + \frac{z_2 - x}{2} \right\|^2 \\ &= 2 \left(\left\| \frac{z_1 - x}{2} \right\|^2 + \left\| \frac{z_2 - x}{2} \right\|^2 \right) - \left\| \frac{z_1 - z_2}{2} \right\|^2 \\ &= \frac{1}{2} \|z_1 - x\|^2 + \frac{1}{2} \|z_2 - x\|^2 - \left\| \frac{z_1 - z_2}{2} \right\|^2 \\ &\leq \|z_1 - x\|^2 - \left\| \frac{z_1 - z_2}{2} \right\|^2, \end{aligned}$$

where the last inequality holds because of 22.5. Hence $z_1 = z_2$ after all and this shows uniqueness.

Now let $\lambda = \inf\{\|x - y\| : y \in K\}$ and let y_n be a minimizing sequence. This means $\{y_n\} \subseteq K$ satisfies $\lim_{n \rightarrow \infty} \|x - y_n\| = \lambda$. By the parallelogram identity,

$$\|y_n - x + y_m - x\|^2 + \|y_n - y_m\|^2 = 2\left(\|y_n - x\|^2 + \|y_m - x\|^2\right)$$

and so, since $\|y_n - x + y_m - x\|^2 = 4\left(\left\|\frac{y_n + y_m}{2} - x\right\|^2\right)$,

$$\begin{aligned} \|y_n - y_m\|^2 &= 2\left(\|y_n - x\|^2 + \|y_m - x\|^2\right) - 4\left(\left\|\frac{y_n + y_m}{2} - x\right\|^2\right) \\ &\leq 2\left(\|y_n - x\|^2 + \|y_m - x\|^2\right) - 4\lambda^2 \end{aligned}$$

The right side converges as $m, n \rightarrow 0$ to 0. Therefore, $\{y_n\}_{n=1}^\infty$ is a Cauchy sequence. Since H is complete, $y_n \rightarrow y$ for some $y \in H$ which must be in K because K is closed. Therefore

$$\|x - y\| = \lim_{n \rightarrow \infty} \|x - y_n\| = \lambda.$$

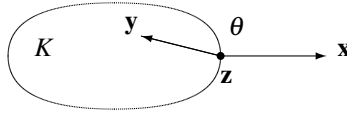
Let $Px = y$. ■

Corollary 22.1.9 *Let K be a closed, convex, nonempty subset of a Hilbert space, H , and let $x \in H$. Then for $z \in K$, $z = Px$ if and only if*

$$\operatorname{Re}(x - z, y - z) \leq 0 \quad (22.6)$$

for all $y \in K$.

Before proving this, consider what it says in the case where the Hilbert space is \mathbb{R}^n .



Condition 22.6 says the angle θ , shown in the diagram, is always obtuse. Remember from calculus, the sign of $x \cdot y$ is the same as the sign of the cosine of the included angle between x and y . Thus, in finite dimensions, the conclusion of this corollary says that $z = Px$ exactly when the angle of the indicated angle is obtuse. Surely the picture suggests this is reasonable.

The inequality 22.6 is an example of a variational inequality and this corollary characterizes the projection of x onto K as the solution of this variational inequality.

Proof of Corollary: Let $z \in K$ and let $y \in K$ also. Since K is convex, it follows that if $t \in [0, 1]$,

$$z + t(y - z) = (1 - t)z + ty \in K.$$

Furthermore, every point of K can be written in this way. (Let $t = 1$ and $y \in K$.) Therefore, $z = Px$ if and only if for all $y \in K$ and $t \in [0, 1]$,

$$\|x - (z + t(y - z))\|^2 = \|(x - z) - t(y - z)\|^2 \geq \|x - z\|^2$$

for all $t \in [0, 1]$ and $y \in K$ if and only if for all $t \in [0, 1]$ and $y \in K$

$$\|x - z\|^2 + t^2\|y - z\|^2 - 2t \operatorname{Re}(x - z, y - z) \geq \|x - z\|^2$$

If and only if for all $t \in [0, 1]$,

$$t^2 \|y - z\|^2 - 2t \operatorname{Re}(x - z, y - z) \geq 0. \quad (22.7)$$

Now this is equivalent to 22.7 holding for all $t \in (0, 1)$. Therefore, dividing by $t \in (0, 1)$, 22.7 is equivalent to

$$t \|y - z\|^2 - 2 \operatorname{Re}(x - z, y - z) \geq 0$$

for all $t \in (0, 1)$ which is equivalent to 22.6. ■

Corollary 22.1.10 *Let K be a nonempty convex closed subset of a Hilbert space, H . Then the projection map P is continuous. In fact, $|Px - Py| \leq |x - y|$.*

Proof: Let $x, x' \in H$. Then by Corollary 22.1.9,

$$\operatorname{Re}(x' - Px', Px - Px') \leq 0, \operatorname{Re}(x - Px, Px' - Px) \leq 0$$

Hence

$$\begin{aligned} 0 &\leq \operatorname{Re}(x - Px, Px - Px') - \operatorname{Re}(x' - Px', Px - Px') \\ &= \operatorname{Re}(x - x', Px - Px') - |Px - Px'|^2 \end{aligned}$$

and so $|Px - Px'|^2 \leq |x - x'| |Px - Px'|$. ■

The next corollary is a more general form for the Brouwer fixed point theorem.

Corollary 22.1.11 *Let $f : K \rightarrow K$ where K is a convex compact subset of \mathbb{R}^n . Then f has a fixed point.*

Proof: Let $K \subseteq \overline{B(0, R)}$ and let P be the projection map onto K . Then consider the map $f \circ P$ which maps $\overline{B(0, R)}$ to $\overline{B(0, R)}$ and is continuous. By the Brouwer fixed point theorem for balls, this map has a fixed point. Thus there exists x such that $(f \circ P)(x) = x$. Now the equation also requires $x \in K$ and so $P(x) = x$. Hence $f(x) = x$. ■

Recall the following definition from linear algebra about direct sum notation.

Definition 22.1.12 *Let H be a vector space and let U and V be subspaces. $U \oplus V = H$ if every element of H can be written as a sum of an element of U and an element of V in a unique way.*

The case where the closed convex set is a closed subspace is of special importance and in this case the above corollary implies the following.

Corollary 22.1.13 *Let K be a closed subspace of a Hilbert space H , and let $x \in H$. Then for $z \in K$, $z = Px$ if and only if*

$$(x - z, y) = 0 \quad (22.8)$$

for all $y \in K$. Furthermore, $H = K \oplus K^\perp$ where

$$K^\perp \equiv \{x \in H : (x, k) = 0 \text{ for all } k \in K\}$$

and

$$\|x\|^2 = \|x - Px\|^2 + \|Px\|^2. \quad (22.9)$$

Proof: Since K is a subspace, the condition 22.6 implies $\operatorname{Re}(x - z, y) \leq 0$ for all $y \in K$. Replacing y with $-y$, it follows $\operatorname{Re}(x - z, -y) \leq 0$ which implies $\operatorname{Re}(x - z, y) \geq 0$ for all y . Therefore, $\operatorname{Re}(x - z, y) = 0$ for all $y \in K$. Now let $|\alpha| = 1$ and $\alpha(x - z, y) = |(x - z, y)|$. Since K is a subspace, it follows $\bar{\alpha}y \in K$ for all $y \in K$. Therefore,

$$0 = \operatorname{Re}(x - z, \bar{\alpha}y) = (x - z, \bar{\alpha}y) = \alpha(x - z, y) = |(x - z, y)|.$$

This shows that $z = Px$, if and only if 22.8.

For $x \in H$, $x = x - Px + Px$ and from what was just shown, $x - Px \in K^\perp$ and $Px \in K$. This shows that $K^\perp + K = H$. Is there only one way to write a given element of H as a sum of a vector in K with a vector in K^\perp ? Suppose $y + z = y_1 + z_1$ where $z, z_1 \in K^\perp$ and $y, y_1 \in K$. Then $(y - y_1) = (z_1 - z)$ and so from what was just shown, $(y - y_1, y - y_1) = (y - y_1, z_1 - z) = 0$ which shows $y_1 = y$ and consequently $z_1 = z$. Finally, letting $z = Px$,

$$\begin{aligned} \|x\|^2 &= (x - z + z, x - z + z) = \|x - z\|^2 + (x - z, z) + (z, x - z) + \|z\|^2 \\ &= \|x - z\|^2 + \|z\|^2 \blacksquare \end{aligned}$$

The following theorem is called the Riesz representation theorem for the dual of a Hilbert space. If $z \in H$ then define an element $f \in H'$ by the rule $(x, z) \equiv f(x)$. It follows from the Cauchy Schwarz inequality and the properties of the inner product that $f \in H'$. The Riesz representation theorem says that all elements of H' are of this form.

Theorem 22.1.14 *Let H be a Hilbert space and let $f \in H'$. Then there exists a unique $z \in H$ such that $f(x) = (x, z)$ for all $x \in H$.*

Proof: Letting $y, w \in H$ the assumption that f is linear implies

$$f(yf(w) - f(y)w) = f(w)f(y) - f(y)f(w) = 0$$

which shows that $yf(w) - f(y)w \in f^{-1}(0)$, which is a closed subspace of H since f is continuous. If $f^{-1}(0) = H$, then f is the zero map and $z = 0$ is the unique element of H which satisfies $f(x) = (x, z)$. If $f^{-1}(0) \neq H$, pick $u \notin f^{-1}(0)$ and let $w \equiv u - Pu \neq 0$. Thus Corollary 22.1.13 implies $(y, w) = 0$ for all $y \in f^{-1}(0)$. In particular, let $y = xf(w) - f(x)w$ where $x \in H$ is arbitrary. Therefore,

$$0 = (f(w)x - f(x)w, w) = f(w)(x, w) - f(x)\|w\|^2.$$

Thus, solving for $f(x)$ and using the properties of the inner product,

$$f(x) = \left(x, \frac{\overline{f(w)}w}{\|w\|^2} \right)$$

Let $z = \overline{f(w)}w/\|w\|^2$. This proves the existence of z . If $f(x) = (x, z_i)$ $i = 1, 2$, for all $x \in H$, then for all $x \in H$, then $(x, z_1 - z_2) = 0$ which implies, upon taking $x = z_1 - z_2$ that $z_1 = z_2$. \blacksquare

If $R: H \rightarrow H'$ is defined by $Rx(y) \equiv (y, x)$, the Riesz representation theorem above states this map is onto. This map is called the Riesz map. It is routine to show R is conjugate linear and $\|Rx\| = \|x\|$. In fact,

$$\begin{aligned} R(\alpha x + \beta y)(u) &\equiv (u, \alpha x + \beta y) = \bar{\alpha}(u, x) + \bar{\beta}(u, y) \\ &\equiv \bar{\alpha}Rx(u) + \bar{\beta}Ry(u) = (\bar{\alpha}Rx + \bar{\beta}Ry)(u) \end{aligned}$$

so it is conjugate linear meaning it goes across plus signs and you factor out conjugates.

$$\|Rx\| \equiv \sup_{\|y\| \leq 1} |Rx(y)| \equiv \sup_{\|y\| \leq 1} |(y, x)| = \|x\|$$

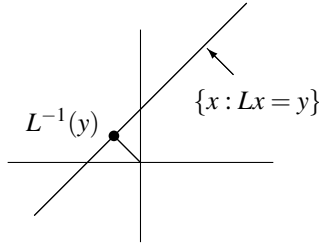
22.2 The Hilbert Space $L(U)$

Let $L \in \mathcal{L}(U, H)$. Then one can consider the image of $L, L(U)$ as a Hilbert space. This is another interesting application of Theorem 22.1.8. First here is a definition which involves abominable and atrociously misleading notation which nevertheless seems to be well accepted.

Definition 22.2.1 Let $L \in \mathcal{L}(U, H)$, the bounded linear maps from U to H for U, H Hilbert spaces. For $y \in L(U)$, let $L^{-1}y$ denote the unique vector in

$$\{x \in U : Lx = y\} \equiv M_y$$

which is closest in U to 0.



Note this is a good definition because $\{x \in U : Lx = y\}$ is closed thanks to the continuity of L and it is obviously convex. Thus Theorem 22.1.8 applies. With this definition define an inner product on $L(U)$ as follows. For $y, z \in L(U)$,

$$(y, z)_{L(U)} \equiv (L^{-1}y, L^{-1}z)_U$$

The notation is abominable because $L^{-1}(y)$ is the normal notation for M_y .

In terms of linear algebra, this L^{-1} is the Moore Penrose inverse. There you obtain the least squares solution x to $Lx = y$ which has smallest norm. Here there is an actual solution and among those solutions you get the one which has least norm. Of course a real honest solution is also a least squares solution so this is the Moore Penrose inverse restricted to $L(U)$.

Lemma 22.2.2 In the context of the above definition, $L^{-1}(y)$ is characterized by

$$\begin{aligned} (L^{-1}(y), x)_U &= 0 \text{ for all } x \in \ker(L) \\ L(L^{-1}(y)) &= y, \quad (L^{-1}(y) \in M_y) \end{aligned}$$

In addition to this, L^{-1} is linear and the above definition does define an inner product.

Proof: By definition, $L^{-1}(y)$ is the unique point of the closed convex set M_y which is closest to 0 in U . Thus it is characterized by

$$(0 - L^{-1}(y), u - L^{-1}(y))_U \leq 0$$

for all $u \in M_y$. Note that $L(u - L^{-1}(y)) = y - y = 0$. Also, if $v \in \ker(L)$, then if $u = L^{-1}(y) + v$, then $u - L^{-1}(y) \in \ker(L)$. Thus a generic element of $\ker(L)$ is $u - L^{-1}(y)$

for $u \in M_y$ and $L^{-1}(y)$ is therefore characterized by $(L^{-1}(y), v)_U = 0$ for all $v \in \ker(L)$ because $\ker(L)$ is a subspace. Also, $L(L^{-1}(y)) = y$. Now from this characterization of L^{-1} , it is obvious that L^{-1} is linear. The inner product is well defined because $L^{-1}(y)$ is uniquely determined. Does it satisfy the axioms? Say $0 = (y, y)_{L(U)}$. Then $L^{-1}(y) = 0$ and so doing L to both sides, $y = 0$. It is clear that $(y, z)_{L(U)} = \overline{(z, y)_{L(U)}}$ because these are defined in terms of a given inner product on U .

$$\begin{aligned} (ay + b\hat{y}, z)_{L(U)} &\equiv (L^{-1}(ay + b\hat{y}), L^{-1}z)_U \\ &= (aL^{-1}(y) + bL^{-1}(\hat{y}), L^{-1}z)_U \\ &= a(L^{-1}(y), L^{-1}z)_U + b(L^{-1}(\hat{y}), L^{-1}z)_U \\ &= a(y, z)_{L(U)} + b(\hat{y}, z)_{L(U)} \end{aligned}$$

Thus this is an inner product as claimed. ■

With the above definition, here is the main result.

Theorem 22.2.3 *Let U, H be Hilbert spaces and let $L \in \mathcal{L}(U, H)$. Then Definition 22.2.1 makes $L(U)$ into a Hilbert space. Also $L : U \rightarrow L(U)$ is continuous and $L^{-1} : L(U) \rightarrow U$ is continuous. Also,*

$$\|L\|_{\mathcal{L}(U, H)} \|Lx\|_{L(U)} \geq \|Lx\|_H \quad (22.10)$$

If U is separable, so is $L(U)$. Also $(L^{-1}(y), x) = 0$ for all $x \in \ker(L)$, and $L^{-1} : L(U) \rightarrow U$ is linear. Also, in case that L is one to one, both L and L^{-1} preserve norms.

Proof: First consider the claim that $L : U \rightarrow L(U)$ is continuous and $L^{-1} : L(U) \rightarrow U$ is also continuous.

$$\|Lu\|_{L(U)}^2 = (L^{-1}(Lu), L^{-1}(Lu))_U \leq \|u\|_U^2$$

(Recall that $L^{-1}(Lu)$ is the smallest vector in U which maps to Lu . Since u is mapped by L to Lu , it follows that $\|L^{-1}(L(u))\|_U \leq \|u\|_U$.) Hence L is continuous.

Next, why is L^{-1} continuous? By definition of the norm, $\|L^{-1}(y)\|_U^2 \equiv \|y\|_{L(U)}^2$. Thus L^{-1} is continuous and $\|L^{-1}\|_{\mathcal{L}(L(U), U)} = 1$.

Why is $L(U)$ a Hilbert space? Let $\{y_n\}$ be a Cauchy sequence in $L(U)$.

$$\|y_n - y_m\|_{L(U)}^2 \equiv \|L^{-1}(y_n - y_m)\|_U^2$$

Then from what was just observed, it follows that $L^{-1}(y_n)$ is a Cauchy sequence in U . Hence $L^{-1}(y_n) \rightarrow x \in U$. Then by continuity of L just shown, $y_n \rightarrow Lx$. This shows that $L(U)$ is a Hilbert space. It was already shown that it is an inner product space and this has shown that it is complete.

If $x \in U$, then $\|Lx\|_H \leq \|L\|_{\mathcal{L}(U, H)} \|x\|_U$. It follows that

$$\begin{aligned} \|L(x)\|_H &= \|L(L^{-1}(L(x)))\|_H \leq \|L\|_{\mathcal{L}(U, H)} \|L^{-1}(L(x))\|_U \\ &= \|L\|_{\mathcal{L}(U, H)} \|L(x)\|_{L(U)}. \end{aligned}$$

This verifies 22.10.

If U is separable, then letting D be a countable dense subset, it follows from the continuity of the operators L, L^{-1} discussed above that $L(D)$ is separable in $L(U)$. To see this, note that

$$\begin{aligned} \|Lx_n - Lx\|_{L(U)} &= \|L(L^{-1}(Lx_n - Lx))\| \leq \|L\|_{\mathcal{L}(U, H)} \|L^{-1}(L(x_n - x))\|_U \\ &\leq \|L\|_{\mathcal{L}(U, H)} \|x_n - x\|_U \end{aligned}$$

As before, $L^{-1}(L(x_n - x))$ is the smallest vector which maps onto $L(x_n - x)$ and so its norm is no larger than $\|x_n - x\|_U$.

Consider the last claim. If L is one to one, then for $y \in L(U)$, there is only one vector which maps to y . Therefore, $L^{-1}(L(x)) = x$. Hence for $y \in L(U)$, $\|y\|_{L(U)} \equiv \|L^{-1}(y)\|_U$. Also,

$$\|Lu\|_{L(U)} \equiv \|L^{-1}(L(u))\|_U \equiv \|u\|_U$$

Thus when L is one to one, $\|L\|_{\mathcal{L}(U, L(U))} = 1$. ■

22.3 Approximations in Hilbert Space

The Gram Schmidt process applies in any vector space which has an inner product.

Theorem 22.3.1 *Let $\{x_1, \dots, x_n\}$ be a basis for M a subspace of H a Hilbert space. Then there exists an orthonormal basis for M , $\{u_1, \dots, u_n\}$ which has the property that for each $k \leq n$, $\text{span}(x_1, \dots, x_k) = \text{span}(u_1, \dots, u_k)$. Also if $\{x_1, \dots, x_n\} \subseteq H$, then the finite dimensional subspace $\text{span}(x_1, \dots, x_n)$ is a closed subspace.*

Proof: Let $\{x_1, \dots, x_n\}$ be a basis for M . Let $u_1 \equiv x_1 / |x_1|$. Thus for $k = 1$, $\text{span}(u_1) = \text{span}(x_1)$ and $\{u_1\}$ is an orthonormal set. Now suppose for some $k < n$, u_1, \dots, u_k have been chosen such that $(u_j \cdot u_l) = \delta_{jl}$ and $\text{span}(x_1, \dots, x_k) = \text{span}(u_1, \dots, u_k)$. Then define

$$u_{k+1} \equiv \frac{x_{k+1} - \sum_{j=1}^k (x_{k+1} \cdot u_j) u_j}{\left| x_{k+1} - \sum_{j=1}^k (x_{k+1} \cdot u_j) u_j \right|}, \quad (22.11)$$

where the denominator is not equal to zero because the x_j form a basis and so

$$x_{k+1} \notin \text{span}(x_1, \dots, x_k) = \text{span}(u_1, \dots, u_k)$$

Thus by induction,

$$u_{k+1} \in \text{span}(u_1, \dots, u_k, x_{k+1}) = \text{span}(x_1, \dots, x_k, x_{k+1}).$$

Also, $x_{k+1} \in \text{span}(u_1, \dots, u_k, u_{k+1})$ which is seen easily by solving 22.11 for x_{k+1} and it follows

$$\text{span}(x_1, \dots, x_k, x_{k+1}) = \text{span}(u_1, \dots, u_k, u_{k+1}).$$

If $l \leq k$,

$$\begin{aligned} (u_{k+1} \cdot u_l) &= C \left((x_{k+1} \cdot u_l) - \sum_{j=1}^k (x_{k+1} \cdot u_j) (u_j \cdot u_l) \right) \\ &= C \left((x_{k+1} \cdot u_l) - \sum_{j=1}^k (x_{k+1} \cdot u_j) \delta_{lj} \right) = C((x_{k+1} \cdot u_l) - (x_{k+1} \cdot u_l)) = 0. \end{aligned}$$

The vectors, $\{u_j\}_{j=1}^n$, generated in this way are therefore an orthonormal basis because each vector has unit length.

Consider the second claim about finite dimensional subspaces. Without loss of generality, assume $\{x_1, \dots, x_n\}$ is linearly independent. If it is not, delete vectors until a linearly independent set is obtained. Then by the first part,

$$\text{span}(x_1, \dots, x_n) = \text{span}(u_1, \dots, u_n) \equiv M$$

where the u_i are an orthonormal set of vectors. Suppose $\{y_k\} \subseteq M$ and $y_k \rightarrow y \in H$. Is $y \in M$? Let $y_k \equiv \sum_{j=1}^n c_j^k u_j$. Then let $\mathbf{c}^k \equiv (c_1^k, \dots, c_n^k)^T$. Then

$$\begin{aligned} \|\mathbf{c}^k - \mathbf{c}^l\|^2 &= \sum_{j=1}^n |c_j^k - c_j^l|^2 = \left(\sum_{j=1}^n (c_j^k - c_j^l) u_j, \sum_{j=1}^n (c_j^k - c_j^l) u_j \right) \\ &= \|y_k - y_l\|^2 \end{aligned}$$

which shows $\{\mathbf{c}^k\}$ is a Cauchy sequence in \mathbb{F}^n and so it converges to $\mathbf{c} \in \mathbb{F}^n$. Thus

$$y = \lim_{k \rightarrow \infty} y_k = \lim_{k \rightarrow \infty} \sum_{j=1}^n c_j^k u_j = \sum_{j=1}^n c_j u_j \in M. \blacksquare$$

Theorem 22.3.2 Let M be the span of $\{u_1, \dots, u_n\}$ in a Hilbert space H and let $y \in H$. Then Py is given by

$$Py = \sum_{k=1}^n (y, u_k) u_k \quad (22.12)$$

and the distance is given by

$$\sqrt{|y|^2 - \sum_{k=1}^n |(y, u_k)|^2}. \quad (22.13)$$

Proof:

$$\begin{aligned} \left(y - \sum_{k=1}^n (y, u_k) u_k, u_p \right) &= (y, u_p) - \sum_{k=1}^n (y, u_k) (u_k, u_p) \\ &= (y, u_p) - (y, u_p) = 0 \end{aligned}$$

It follows that $(y - \sum_{k=1}^n (y, u_k) u_k, u) = 0$ for all $u \in M$ and so by Corollary 22.1.13 this verifies 22.12.

The square of the distance, d is given by

$$\begin{aligned} d^2 &= \left(y - \sum_{k=1}^n (y, u_k) u_k, y - \sum_{k=1}^n (y, u_k) u_k \right) \\ &= |y|^2 - 2 \sum_{k=1}^n |(y, u_k)|^2 + \sum_{k=1}^n |(y, u_k)|^2 \end{aligned}$$

and this shows 22.13. \blacksquare

22.4 Orthonormal Sets

The concept of an orthonormal set of vectors is a generalization of the notion of the standard basis vectors of \mathbb{R}^n or \mathbb{C}^n .

Definition 22.4.1 Let H be a Hilbert space. $S \subseteq H$ is called an orthonormal set if $\|x\| = 1$ for all $x \in S$ and $(x, y) = 0$ if $x, y \in S$ and $x \neq y$. For any set, D ,

$$D^\perp \equiv \{x \in H : (x, d) = 0 \text{ for all } d \in D\}.$$

If S is a set, $\text{span}(S)$ is the set of all finite linear combinations of vectors from S .

You should verify that D^\perp is always a closed subspace of H . It is assumed that our Hilbert spaces are not $\{0\}$.

Theorem 22.4.2 In any separable Hilbert space H , there exists a countable orthonormal set, $S = \{x_i\}$ such that the span of these vectors is dense in H . Furthermore, if $\text{span}(S)$ is dense, then for $x \in H$,

$$x = \sum_{i=1}^{\infty} (x, x_i) x_i \equiv \lim_{n \rightarrow \infty} \sum_{i=1}^n (x, x_i) x_i. \quad (22.14)$$

$$\text{Also, } (x, y) = \sum_{i=1}^{\infty} (x, x_i) \overline{(y, x_i)}.$$

Proof: Let \mathcal{F} denote the collection of all orthonormal subsets of H . \mathcal{F} is nonempty because some $\{x\} \in \mathcal{F}$ where $\|x\| = 1$. The set, \mathcal{F} is a partially ordered set with the order given by set inclusion. By the Hausdorff maximal theorem, there exists a maximal chain, \mathcal{C} in \mathcal{F} . Then let $S \equiv \bigcup \mathcal{C}$. It follows S must be a maximal orthonormal set of vectors. This is because if $x, y \in S$, then both vectors are in a single $C \in \mathcal{C}$. Therefore, $(x, y) = 0$ or one.

It remains to verify that S is countable $\text{span}(S)$ is dense, and the condition, 22.14 holds. To see S is countable note that if $x, y \in S$, then

$$\|x - y\|^2 = \|x\|^2 + \|y\|^2 - 2\text{Re}(x, y) = \|x\|^2 + \|y\|^2 = 2.$$

Therefore, the open sets, $B(x, \frac{1}{2})$ for $x \in S$ are disjoint and cover S . Since H is assumed to be separable, there exists a point from a countable dense set in each of these disjoint balls showing there can only be countably many of the balls and that consequently, S is countable as claimed.

It remains to verify 22.14 and that $\text{span}(S)$ is dense. If $\text{span}(S)$ is not dense, then $\overline{\text{span}(S)}$ is a closed proper subspace of H and letting $y \notin \overline{\text{span}(S)}$, $z \equiv \frac{y - Py}{\|y - Py\|} \in \text{span}(S)^\perp$ where P is the projection map mentioned earlier. But then $S \cup \{z\}$ would be a larger orthonormal set of vectors contradicting the maximality of S .

It remains to verify 22.14. Let $S = \{x_i\}_{i=1}^{\infty}$ and consider the problem of choosing the constants, c_k in such a way as to minimize the expression

$$\left\| x - \sum_{k=1}^n c_k x_k \right\|^2 = \|x\|^2 + \sum_{k=1}^n |c_k|^2 - \sum_{k=1}^n \overline{c_k} (x, x_k) - \sum_{k=1}^n c_k \overline{(x, x_k)}.$$

This equals $\|x\|^2 + \sum_{k=1}^n |c_k - (x, x_k)|^2 - \sum_{k=1}^n |(x, x_k)|^2$ and therefore, this minimum is achieved when $c_k = (x, x_k)$ and equals $\|x\|^2 - \sum_{k=1}^n |(x, x_k)|^2$. Now since $\text{span}(S)$ is dense, if n large enough then for some choice of constants c_k , $\|x - \sum_{k=1}^n c_k x_k\|^2 < \varepsilon$. However, from what was just shown,

$$\left\| x - \sum_{i=1}^n (x, x_i) x_i \right\|^2 \leq \left\| x - \sum_{k=1}^n c_k x_k \right\|^2 < \varepsilon$$

showing that $\lim_{n \rightarrow \infty} \sum_{i=1}^n (x, x_i) x_i = x$ as claimed.

For the last claim, from what was just shown and the observation that if $x_n \rightarrow x$ and $y_n \rightarrow y$, then $(x_n, y_n) \rightarrow (x, y)$,

$$\begin{aligned} (x, y) &= \lim_{n \rightarrow \infty} \left(\sum_{i=1}^n (x, x_i) x_i, \sum_{j=1}^n (y, x_j) x_j \right) = \lim_{n \rightarrow \infty} \sum_{i, j \leq n} (x, x_i) \overline{(y, x_j)} (x_i, x_j) \\ &= \lim_{n \rightarrow \infty} \sum_{i=1}^n (x, x_i) \overline{(y, x_i)} = \sum_{i=1}^{\infty} (x, x_i) \overline{(y, y_i)} \blacksquare \end{aligned}$$

The proof of this theorem contains the following corollary.

Corollary 22.4.3 *Let S be any orthonormal set of vectors and let*

$$\{x_1, \dots, x_n\} \subseteq S.$$

Then if $x \in H$

$$\left\| x - \sum_{k=1}^n c_k x_k \right\|^2 \geq \left\| x - \sum_{i=1}^n (x, x_i) x_i \right\|^2$$

for all choices of constants, c_k . In addition to this, Bessel's inequality follows,

$$\|x\|^2 \geq \sum_{k=1}^n |(x, x_k)|^2.$$

If S is countable and $\text{span}(S)$ is dense, then letting $\{x_i\}_{i=1}^{\infty} = S$, 22.14 follows.

Corollary 22.4.4 *If V is a closed subspace of a Hilbert space H and if*

$$V = \overline{\text{span}(\{u_k\}_{k=1}^N)}$$

where the u_k are an orthonormal set and $N \leq \infty$, then for P the projection map onto V , it follows that $Py = \sum_{k=1}^N (y, u_k) u_k$ and when $N = \infty$ the series converges to Py . In particular, when $y \in V$, $y = Py = \sum_{k=1}^N (y, u_k) u_k$.

Proof: The case where $N < \infty$ was done above. For $y \in H$,

$$\left(y - \sum_{k=1}^{\infty} (y, u_k) u_k, u_m \right) = (y, u_m) - (y, u_m) = 0$$

assuming the series makes sense. Therefore, if this happens, then $Py = \sum_{k=1}^{\infty} (y, u_k) u_k$. Now note that $(y, u_k) = (y - Py, u_k) + (Py, u_k) = (Py, u_k)$. Now, by assumption, there are scalars

c_k^n such that $\lim_{n \rightarrow \infty} \|Py - \sum_{k=1}^n c_k^n u_k\| = 0$. Then by Corollary 22.4.3 and what was just observed,

$$0 = \lim_{n \rightarrow \infty} \left\| Py - \sum_{k=1}^n (Py, u_k) u_k \right\| = \lim_{n \rightarrow \infty} \left\| Py - \sum_{k=1}^n (y, u_k) u_k \right\|$$

and so the sum converges to Py as claimed. ■

22.5 Compact Operators in Hilbert Space

Definition 22.5.1 Let $A \in \mathcal{L}(H, H)$ where H is a Hilbert space. So the map, $x \rightarrow (Ax, y)$ is continuous and linear. By the Riesz representation theorem, there exists a unique element of H , denoted by A^*y such that

$$(Ax, y) = (x, A^*y).$$

It is clear $y \rightarrow A^*y$ is linear and continuous. It is linear because

$$\begin{pmatrix} (x, A^*(ay + bz)) = \bar{a}(Ax, y) + \bar{b}(Ax, z) \\ = \bar{a}(x, A^*y) + \bar{b}(x, A^*z) = (x, aA^*y + bA^*z) \end{pmatrix}$$

A^* is called the adjoint of A . A is a self adjoint operator if $A = A^*$. Thus for a self adjoint operator, $(Ax, y) = (x, Ay)$ for all $x, y \in H$ and so $(Ax, x) = \overline{(x, Ax)} = \overline{(Ax, x)}$ so (Ax, x) is real. A is a compact operator if whenever $\{x_k\}$ is a bounded sequence, there exists a convergent subsequence of $\{Ax_k\}$. Equivalently, A maps bounded sets to sets whose closures are compact.

There is an important observation about the range of a compact operator. It is a general result so I will express it in terms of Banach space.

Proposition 22.5.2 Let X be a Banach space and let $A \in \mathcal{L}(X, X)$ be a compact operator meaning that the image of a bounded set is a pre-compact set. Then $A(X)$ is separable.

Proof: If for every $n \in \mathbb{N}$ there is a $1/n$ net for $A(B(0, 1))$, then $A(B(0, 1))$ would be separable, and a countable dense subset would be the union of these $1/n$ nets. It would follow then that $A(X)$ is also separable. A countable dense subset would be positive rational numbers times these countably many points in $A(B(0, 1))$. Therefore, if $A(X)$ is not separable, there is some $\varepsilon > 0$ such that there is no ε net, for $A(B(0, 1))$ meaning that there are infinitely many points which are all ε apart in $A(B(0, 1))$. Let these points be $\{Au_k\}_{k=1}^\infty$. But now this is a contradiction because there can be no convergent subsequence. ■

The big result is called the Hilbert Schmidt theorem. It is a generalization to arbitrary Hilbert spaces of standard finite dimensional results having to do with diagonalizing a symmetric matrix.

Theorem 22.5.3 Let A be a compact self adjoint operator defined on a Hilbert space H . Then there exists a countable set of eigenvalues, $\{\lambda_i\}$ and an orthonormal set of eigenvectors, u_i satisfying

$$\lambda_i \text{ is real, } |\lambda_n| \geq |\lambda_{n+1}|, Au_i = \lambda_i u_i, \quad (22.15)$$

and either $\lim_{n \rightarrow \infty} \lambda_n = 0$, or for some n , $\text{span}(u_1, \dots, u_n) = A(H)$. In any case,

$$\text{span}(\{u_i\}_{i=1}^\infty) \text{ is dense in } A(H). \quad (22.16)$$

and for all $x \in H$, $Ax = \sum_{k=1}^\infty \lambda_k (x, u_k) u_k$ where the sum might be finite.

Proof: First note that if you have a self adjoint operator A , then its eigenvalues are real. Here is why:

$$(Au, u) = \lambda \|u\|^2 = (u, Au) = \bar{\lambda} \|u\|^2.$$

If $\|A\| = 0$ then pick $u \in H$ with $\|u\| = 1$ and let $\lambda = 0$. Since $A(H) = 0$ it follows the span of u is dense in $A(H)$ and the formula for Ax holds.

Assume $A \neq 0$. I will show there exists an eigenvector u . From the definition of $\|A\|$ there exists $x_n, \|x_n\| = 1$, and $\|Ax_n\| \rightarrow \|A\| \equiv |\lambda|$. Now it is clear that A^2 is also a compact self adjoint operator. Consider

$$\left((\lambda^2 - A^2)x_n, x_n \right) = \lambda^2 (x_n, x_n) - (A^2 x_n, x_n) = \lambda^2 - \|Ax_n\|^2 \rightarrow 0.$$

Since A is compact, there exists a subsequence of $\{x_n\}$ still denoted by $\{x_n\}$ such that Ax_n converges to some element of H . Thus since $\lambda^2 - A^2$ satisfies $\left((\lambda^2 - A^2)y, y \right) \geq 0$ in addition to being self adjoint, it follows

$$x, y \rightarrow \left((\lambda^2 - A^2)x, y \right)$$

satisfies all the axioms for an inner product except for the one which says that $(z, z) = 0$ only if $z = 0$. Therefore, the Cauchy Schwarz inequality may be used to write

$$\left| \left((\lambda^2 - A^2)x_n, y \right) \right| \leq \left((\lambda^2 - A^2)y, y \right)^{1/2} \left((\lambda^2 - A^2)x_n, x_n \right)^{1/2} \leq e_n \|y\|.$$

where $e_n \rightarrow 0$ as $n \rightarrow \infty$. Taking the sup over $\|y\| \leq 1$, $\lim_{n \rightarrow \infty} \left\| (\lambda^2 - A^2)x_n \right\| = 0$. Since $A^2 x_n$ converges, it follows, since $\lambda \neq 0$ that $\{x_n\}$ is a Cauchy sequence converging to x with $\|x\| = 1$. Therefore, $A^2 x_n \rightarrow A^2 x$ and so $\left\| (\lambda^2 - A^2)x \right\| = 0$. Now this shows that

$$(\lambda I + A)(\lambda I - A)x = 0.$$

If $(\lambda I - A)x = 0$, let $u \equiv x$. If $(\lambda I - A)x = y \neq 0$, let $u \equiv \frac{y}{\|y\|}$. Note that this did not identify the sign of λ . Also note that since $\lambda \neq 0, u \in A(H)$.

Let $\mathcal{A} \in \mathcal{F}$ mean that \mathcal{A} consists of vectors of $A(H)$, \mathcal{A} is an orthonormal set of vectors, and for each $u \in \mathcal{A}, Au = \lambda u$ for some λ . I claim that \mathcal{A} is countable because from the compactness of $A, A(H)$ is separable by Proposition 22.5.2 but these vectors of $A(H)$ are all at least $1/2$ apart. Partially order \mathcal{F} by set inclusion. Let \mathcal{C} be a maximal chain. Then $\mathcal{A}_\infty \equiv \cup \mathcal{C}$ is a maximal element of \mathcal{F} . I need to show its span is dense in $A(H)$. If $\overline{\text{span } \mathcal{A}_\infty}$ fails to contain $A(H)$, then there is a nonzero vector $w \equiv Av$ which is not in $\overline{\text{span } \mathcal{A}_\infty}$. Then $(w - Pw) / \|w - Pw\|$ is a unit vector perpendicular to $\overline{\text{span } \mathcal{A}_\infty}$. Is this vector in $A(H)$? Is $P(Av) = A(Pv)$? Using Corollary 22.4.4,

$$\begin{aligned} P(w) &= P(Av) = \sum_{k=1}^N (Av, u_k) u_k = \sum_{k=1}^N \lambda_k (v, u_k) u_k \\ &= \sum_{k=1}^N (v, u_k) Au_k = A(Pv) \in A(H) \end{aligned}$$

so this unit vector w is in $(\overline{\text{span } \mathcal{A}_\infty})^\perp$ and $\mathcal{A}_\infty \cup \{w\}$ is larger than the maximal element of \mathcal{F} so it must be the case that $\overline{\text{span } \mathcal{A}_\infty} \supseteq A(H)$ after all.

As noted, this orthonormal set \mathcal{A}_∞ is countable. Let it be $\{u_k\}_{k=1}^N$ where $N \leq \infty$. Thus for $x \in H$, $Ax \in A(H) \subseteq \overline{A(H)}$ and so, by Corollary 22.4.4,

$$Ax = \sum_{k=1}^N (Ax, u_k) u_k = \sum_{k=1}^N (x, Au_k) u_k = \sum_{k=1}^N \lambda_k (x, u_k) u_k \quad (22.17)$$

and the series converges. Also, the formula implies directly that $Au_m = \lambda_m u_m$ so $|\lambda_m| \leq \|A\|$.

I claim $\limsup_{n \rightarrow \infty} |\lambda_n| = 0$. If this were not so, then for some $\varepsilon > 0$, $0 < \varepsilon \leq |\lambda_n|$ for a subsequence still denoted as λ_n but then

$$\|Au_n - Au_m\|^2 = \|\lambda_n u_n - \lambda_m u_m\|^2 = |\lambda_n|^2 + |\lambda_m|^2 \geq 2\varepsilon^2$$

and so there could not exist a convergent subsequence of $\{Au_k\}_{k=1}^\infty$ contrary to the assumption that A is compact. This verifies the claim that $\lim_{n \rightarrow \infty} \lambda_n = 0$. Also, since $|\lambda_m| \leq \|A\|$, if $S \subseteq \mathbb{N}$, $\sup_{m \in S} |\lambda_m| = \max_{m \in S} |\lambda_m|$ and so, we can re-number the u_k if necessary such that the eigenvalues satisfy $|\lambda_k| \geq |\lambda_{k+1}|$ for all k . Thus, if $\lambda_m = 0$ for some m , it follows from 22.17 that $A(H)$ is contained in the span of finitely many of the vectors $\{u_k\}$. ■

Define $v \otimes u \in \mathcal{L}(H, H)$ by $v \otimes u(x) = (x, u)v$, then 22.17 is of the form

$$A = \sum_{k=1}^N \lambda_k u_k \otimes u_k$$

This is the content of the following corollary.

Corollary 22.5.4 *The main conclusion of the above theorem can be written as $A = \sum_{k=1}^N \lambda_k u_k \otimes u_k$ where the convergence of the partial sums takes place in the operator norm.*

Proof: Without loss of generality, assume $N = \infty$.

$$\begin{aligned} & \left| \left(\left(A - \sum_{k=1}^n \lambda_k u_k \otimes u_k \right) x, y \right) \right| = \left| \left(Ax - \sum_{k=1}^n \lambda_k (x, u_k) u_k, y \right) \right| \\ &= \left| \left(\sum_{k=n}^{\infty} \lambda_k (x, u_k) u_k, y \right) \right| = \left| \sum_{k=n}^{\infty} \lambda_k (x, u_k) (u_k, y) \right| \\ &\leq |\lambda_n| \left(\sum_{k=n}^{\infty} |(x, u_k)|^2 \right)^{1/2} \left(\sum_{k=n}^{\infty} |(y, u_k)|^2 \right)^{1/2} \leq |\lambda_n| \|x\| \|y\| \end{aligned}$$

It follows $\|(A - \sum_{k=1}^n \lambda_k u_k \otimes u_k)(x)\| \leq |\lambda_n| \|x\|$ ■

Corollary 22.5.5 *Let A be a compact self adjoint operator and*

$$A = \sum_{k=1}^N \lambda_k u_k \otimes u_k$$

where $Au_k = \lambda_k u_k$ with $\|u_k\| = 1$, the $|\lambda_k|$ decreasing. Then $|\lambda_{k+1}| = \|A_k\|$ where A_k is the restriction of A to $\{u_1, \dots, u_k\}^\perp$.

Proof: First note that if $V_k \equiv \{u_1, \dots, u_k\}^\perp$, then $A : V_k \rightarrow V_k$. Thus

$$\|A_k\| \equiv \sup \{\|A_k u\|, \|u\| \leq 1, (u, u_i) = 0, i \leq k\}$$

Since the $|\lambda_k|$ are decreasing, this will be maximized by picking $u = u_{k+1}$ and then the result is just $|\lambda_k|$. ■

Lemma 22.5.6 *If V_λ is the eigenspace for $\lambda \neq 0$ and $B : V_\lambda \rightarrow V_\lambda$ is a compact self adjoint operator with $Bx = \lambda x$ for all $x \in V_\lambda$ then V_λ must be finite dimensional.*

Proof: First note that $B(V_\lambda) = V_\lambda$ because if $u \in V_\lambda$, then $Bu = \lambda u$ and so $B(\frac{u}{\lambda}) = u$. From Theorem 22.5.3, \mathcal{A}_∞ that maximal set might be finite in which case it would yield a finite orthonormal basis for $V_\lambda = B(V_\lambda)$. But it can't be infinite because there is only one eigenvalue and it is not zero so cannot converge to 0. ■

Next is the case of most interest when H is separable. In this case, the eigenfunctions actually give an orthonormal basis for H .

Corollary 22.5.7 *Let A be a compact self adjoint operator defined on a separable Hilbert space, H . Then there exists a countable set of eigenvalues, $\{\lambda_i\}$ and an orthonormal set of eigenvectors $\{v_i\}$ satisfying $Av_i = \lambda_i v_i$, $\|v_i\| = 1$, $\text{span}(\{v_i\}_{i=1}^\infty)$ is dense in H . Furthermore, if $\lambda_i \neq 0$, the space, $V_{\lambda_i} \equiv \{x \in H : Ax = \lambda_i x\}$ is finite dimensional.*

Proof: Let B be the restriction of A to V_{λ_i} . Thus B is a compact self adjoint operator which maps V_λ to V_λ and has only one eigenvalue λ_i on V_{λ_i} . By Lemma 22.5.6, V_λ is finite dimensional. As to the density of some $\text{span}(\{v_i\}_{i=1}^\infty)$ in H , let $W \equiv \overline{\text{span}(\{u_i\})}^\perp$ where $A = \sum_{k=1}^N \lambda_k u_k \otimes u_k$. By Theorem 22.4.2, there is a maximal orthonormal set of vectors, $\{w_i\}_{i=1}^\infty$ whose span is dense in W . There are only countably many of these since the space H is separable. Then consider $\{v_i\}_{i=1}^\infty = \{u_i\}_{i=1}^\infty \cup \{w_i\}_{i=1}^\infty$. $Aw_i = \sum_{k=1}^N \lambda_k (w_i, u_k) u_k = 0$. Thus each w_i is an eigenvector for A . ■

Suppose $\lambda \notin \{\lambda_k\}_{k=1}^\infty$, the eigenvalues of A , and $\lambda \neq 0$. Then the above formula for A , yields an interesting formula for $(A - \lambda I)^{-1}$. Note first that since $\lim_{n \rightarrow \infty} \lambda_n = 0$, it follows that $\lambda_n^2 / (\lambda_n - \lambda)^2$ must be bounded, say by a positive constant, M .

Corollary 22.5.8 *Let A be a compact self adjoint operator and let $\lambda \notin \{\lambda_n\}_{n=1}^\infty$ and $\lambda \neq 0$ where the λ_n are the eigenvalues of A . ($Ax = \lambda x, x \neq 0$) Then*

$$(A - \lambda I)^{-1} x = -\frac{1}{\lambda} x + \frac{1}{\lambda} \sum_{k=1}^{\infty} \frac{\lambda_k}{\lambda_k - \lambda} (x, u_k) u_k. \quad (22.18)$$

Proof: Let $m < n$. Then since the $\{u_k\}$ form an orthonormal set,

$$\left| \sum_{k=m}^n \frac{\lambda_k}{\lambda_k - \lambda} (x, u_k) u_k \right| \leq \left(\sum_{k=m}^n \left(\frac{\lambda_k}{\lambda_k - \lambda} \right)^2 |(x, u_k)|^2 \right)^{1/2} \leq M \left(\sum_{k=m}^n |(x, u_k)|^2 \right)^{1/2}. \quad (22.19)$$

But from Bessel's inequality, $\sum_{k=1}^\infty |(x, u_k)|^2 \leq \|x\|^2$ and so for m large enough, the first term in 22.19 is smaller than ε . This shows the infinite series in 22.18 converges. It is now routine to verify that the formula in 22.18 is the inverse. ■

22.5.1 Nuclear Operators

A very useful idea in linear algebra is the trace of a matrix. It is one of the principle invariants and is just the sum of the entries along the main diagonal. Thus if $A \in \mathcal{L}(\mathbb{C}^n, \mathbb{C}^n)$, the trace is $\sum_i e_i^T A e_i$ where the e_i are the standard orthonormal basis vectors, e_i having a 1 in the i^{th} position and 0 elsewhere. If you used another orthonormal basis, $\{v_i\}$ the trace would be the same because the mapping $v_i \rightarrow e_i$ will preserve lengths and is therefore a unitary transformation. The two computations would involve a similarity transformation. In infinite dimensions when you have a separable Hilbert space, the notion of trace might not make sense. The nuclear operators are those for which it will make sense.

Definition 22.5.9 A self adjoint operator $A \in \mathcal{L}(H, H)$ for H a separable Hilbert space is called a nuclear operator if for some complete orthonormal set, $\{e_k\}$, it follows that $\sum_{k=1}^{\infty} |(Ae_k, e_k)| < \infty$.

We specialize to self adjoint operators because this will ensure that (Ax, x) is real. To begin with here is an interesting lemma.

Lemma 22.5.10 Suppose $\{A_n\}$ is a sequence of compact operators in $\mathcal{L}(X, Y)$ for two Banach spaces, X and Y and suppose $A \in \mathcal{L}(X, Y)$ and $\lim_{n \rightarrow \infty} \|A - A_n\| = 0$. Then A is also compact.

Proof: Let D be a bounded set in X such that $\|b\| \leq C$ for all $b \in D$. I need to verify $A(D)$ is totally bounded. Suppose then it is not. Then there exists $\varepsilon > 0$ and an infinite sequence, $\{Ab_i\}$ where $b_i \in D$ and $\|Ab_i - Ab_j\| \geq \varepsilon$ whenever $i \neq j$. Then let n be large enough that $\|A - A_n\| \leq \frac{\varepsilon}{4C}$. Then

$$\begin{aligned} \|A_n b_i - A_n b_j\| &= \|Ab_i - Ab_j + (A_n - A)b_i - (A_n - A)b_j\| \\ &\geq \|Ab_i - Ab_j\| - (\|(A_n - A)b_i\| + \|(A_n - A)b_j\|) \\ &\geq \|Ab_i - Ab_j\| - 2\frac{\varepsilon}{4C}C \geq \frac{\varepsilon}{2}, \end{aligned}$$

a contradiction to A_n being compact. ■

Then one can prove the following lemma. In this lemma, $A \geq 0$ will mean $(Ax, x) \geq 0$.

Lemma 22.5.11 Let $A \geq 0$ be a nuclear operator defined on a separable Hilbert space H . Then A is compact and also, whenever $\{e_k\}$ is a complete orthonormal set,

$$A = \sum_{j=1}^{\infty} \sum_{i=1}^{\infty} (Ae_i, e_j) e_i \otimes e_j.$$

Proof: First consider the formula. Since A is given to be continuous,

$$Ax = A \left(\sum_{j=1}^{\infty} (x, e_j) e_j \right) = \sum_{j=1}^{\infty} (x, e_j) Ae_j,$$

the series converging because $x = \sum_{j=1}^{\infty} (x, e_j) e_j$. Then also since A is self adjoint,

$$\sum_{j=1}^{\infty} \sum_{i=1}^{\infty} (Ae_i, e_j) e_i \otimes e_j (x) \equiv \sum_{j=1}^{\infty} \sum_{i=1}^{\infty} (Ae_i, e_j) (x, e_j) e_i = \sum_{j=1}^{\infty} (x, e_j) \sum_{i=1}^{\infty} (Ae_i, e_j) e_i$$

$$= \sum_{j=1}^{\infty} (x, e_j) \sum_{i=1}^{\infty} (Ae_j, e_i) e_i = \sum_{j=1}^{\infty} (x, e_j) Ae_j$$

Next consider the claim that A is compact. Let $C_A \equiv \left(\sum_{j=1}^{\infty} |(Ae_j, e_j)| \right)^{1/2}$. Let A_n be defined by $A_n \equiv \sum_{j=1}^{\infty} \sum_{i=1}^n (Ae_i, e_j) (e_i \otimes e_j)$. Then A_n has values in $\text{span}(e_1, \dots, e_n)$ and so it must be a compact operator because bounded sets in a finite dimensional space must be precompact. Then

$$\begin{aligned} |(Ax - A_n x, y)| &= \left| \sum_{j=1}^{\infty} \sum_{i=n+1}^{\infty} (Ae_i, e_j) (y, e_j) (e_i, x) \right| = \left| \sum_{j=1}^{\infty} (y, e_j) \sum_{i=n+1}^{\infty} (Ae_i, e_j) (e_i, x) \right| \\ &\leq \left| \sum_{j=1}^{\infty} |(y, e_j)| (Ae_j, e_j)^{1/2} \sum_{i=n+1}^{\infty} (Ae_i, e_i)^{1/2} |(e_i, x)| \right| \\ &\leq \left(\sum_{j=1}^{\infty} |(y, e_j)|^2 \right)^{1/2} \left(\sum_{j=1}^{\infty} |(Ae_j, e_j)| \right)^{1/2} \cdot \left(\sum_{i=n+1}^{\infty} |(x, e_i)|^2 \right)^{1/2} \left(\sum_{i=n+1}^{\infty} |(Ae_i, e_i)| \right)^{1/2} \\ &\leq \|y\| \|x\| C_A \left(\sum_{i=n+1}^{\infty} |(Ae_i, e_i)| \right)^{1/2} \end{aligned}$$

and this shows that if n is sufficiently large, it follows that $|((A - A_n)x, y)| \leq \varepsilon \|x\| \|y\|$ so for such n , $\|A - A_n\| < \varepsilon$. Therefore, $\lim_{n \rightarrow \infty} \|A - A_n\| = 0$ and so A is the limit in operator norm of finite rank bounded linear operators, each of which is compact. Therefore, A is also compact. ■

Definition 22.5.12 The trace of a nuclear operator $A \in \mathcal{L}(H, H)$ such that $A \geq 0$ is defined to equal $\sum_{k=1}^{\infty} (Ae_k, e_k)$ where $\{e_k\}$ is an orthonormal basis for the separable Hilbert space, H .

Theorem 22.5.13 Definition 22.5.12 is well defined and equals $\sum_{j=1}^{\infty} \lambda_j$ where the λ_j are the nonnegative eigenvalues of A .

Proof: Suppose $\{u_k\}$ be the basis of eigenvectors of the Hilbert Schmidt theorem Then $e_k = \sum_{j=1}^{\infty} u_j (e_k, u_j)$. By Lemma 22.5.11 A is compact and so by the Hilbert Schmidt theorem, Theorem 22.5.3, $A = \sum_{k=1}^{\infty} \lambda_k u_k \otimes u_k$ where the u_k are the orthonormal eigenvectors of A which form a complete orthonormal set. Then

$$\begin{aligned} \sum_{k=1}^{\infty} (Ae_k, e_k) &= \sum_{k=1}^{\infty} \left(A(e_k), \sum_{i=1}^{\infty} u_i (e_k, u_i) \right) = \sum_{k=1}^{\infty} \sum_{i=1}^{\infty} (A(e_k), u_i) (e_k, u_i) \\ &= \sum_{k=1}^{\infty} \sum_{i=1}^{\infty} (e_k, Au_i) (e_k, u_i) = \sum_{k=1}^{\infty} \sum_{i=1}^{\infty} \left(e_k, \sum_{j=1}^{\infty} (Au_i, u_j) u_j \right) (e_k, u_i) \\ &= \sum_{k=1}^{\infty} \sum_i \sum_j (Au_i, u_j) (e_k, u_j) (u_i, e_k) = \sum_{k=1}^{\infty} \sum_i \sum_j \lambda_j \delta_{ij} (e_k, u_j) (u_i, e_k) \\ &= \sum_{k=1}^{\infty} \sum_{i=1}^{\infty} \lambda_i |(u_i, e_k)|^2 = \sum_{i=1}^{\infty} \lambda_i \sum_{k=1}^{\infty} |(u_i, e_k)|^2 = \sum_{i=1}^{\infty} \lambda_i \|u_i\|^2 = \sum_{i=1}^{\infty} \lambda_i \quad \blacksquare \end{aligned}$$

This is just like it is for a matrix. Recall the trace of a matrix is the sum of the eigenvalues.

It is also easy to see that in any Hilbert space, there exist nuclear operators. Let $\sum_{k=1}^{\infty} |\lambda_k| < \infty$. Then let $\{e_k\}$ be an orthonormal set of vectors. Let $A \equiv \sum_{k=1}^{\infty} \lambda_k e_k \otimes e_k$, λ_k real. It is not too hard to verify this works.

Much more can be said about nuclear operators.

22.5.2 Hilbert Schmidt Operators

In all of this, H, G will be separable Hilbert spaces.

Definition 22.5.14 Let T be a continuous linear mapping from H to G and whenever $\{e_k\}$ is an orthonormal basis for H , then $\sum_k \|Te_k\|^2 < \infty$. Such an operator is called a Hilbert Schmidt operator. We write $T \in \mathcal{L}_2(H, G)$. Picking an orthonormal basis (complete orthonormal set), define $\|T\|_{\mathcal{L}_2}^2 \equiv \sum_k \|Te_k\|_G^2$.

It is necessary to show that this is well defined and does not depend on the orthonormal basis. For now let the orthonormal basis be fixed.

Lemma 22.5.15 If T is Hilbert Schmidt, then $\|T\|_{\mathcal{L}(H, G)} \leq \|T\|_{\mathcal{L}_2(H, G)}$. Also T^* is Hilbert Schmidt and T^*T is Hilbert Schmidt. In fact, if $T \in \mathcal{L}_2(H, G)$ and $S \in \mathcal{L}(G, H)$ then $ST \in \mathcal{L}_2(H, H)$.

Proof: Pick an orthonormal basis for H , $\{e_k\}$ $\sum_k \|Te_k\|^2 < \infty$ and an orthonormal basis for G , $\{f_k\}$. Then let $x = \sum_{k=1}^n (x, e_k) e_k \equiv \sum_{k=1}^n x_k e_k$. Then

$$\begin{aligned} \|Tx\| &= \left(\sum_{k=1}^{\infty} |(Tx, f_k)|^2 \right)^{1/2} = \left(\sum_{k=1}^{\infty} \left| \left(\sum_{j=1}^n x_j Te_j, f_k \right) \right|^2 \right)^{1/2} \\ &= \left(\sum_{k=1}^{\infty} \left| \sum_{j=1}^n (x_j Te_j, f_k) \right|^2 \right)^{1/2} \leq \sum_{j=1}^n \left(\sum_{k=1}^{\infty} |(x_j Te_j, f_k)|^2 \right)^{1/2} \\ &\leq \sum_{j=1}^n |x_j| \left(\sum_{k=1}^{\infty} |(Te_j, f_k)|^2 \right)^{1/2} \\ &= \sum_{j=1}^n |x_j| \|Te_j\| \leq \left(\sum_{j=1}^n |x_j|^2 \right)^{1/2} \left(\sum_{j=1}^n \|Te_j\|^2 \right)^{1/2} = \|x\| \|T\|_{\mathcal{L}_2} \end{aligned}$$

Therefore, since finite sums of the form $\sum_{k=1}^n x_k e_k$ are dense in H , it follows $\|T\| \leq \|T\|_{\mathcal{L}_2}$.

Letting $\{f_i\}$ be orthonormal in G , $\|Te_k\|^2 = \sum_j |(Te_k, f_j)|^2$ and so

$$\sum_k \|Te_k\|^2 = \sum_k \sum_j |(Te_k, f_j)|^2 = \sum_j \sum_k |(e_k, T^* f_j)|^2 = \sum_j \|T^* f_j\|^2$$

so T^* is also Hilbert Schmidt.

Let $\{e_i\}$ be an orthonormal basis for H so $\sum_i \|Te_i\|^2 < \infty$. Then

$$\sum_i \|T^*Te_i\|^2 \leq \|T^*\|^2 \sum_i \|Te_i\|^2 < \infty$$

For the last claim, let $\{e_k\}$ be an orthonormal basis. Then

$$\sum_k \|ST(e_k)\|^2 \leq \|S\|^2 \sum_k \|Te_k\|^2 < \infty. \blacksquare$$

Definition 22.5.16 Define $(S, T) \equiv \sum_k (Se_k, Te_k)$ where $\{e_k\}$ is a given orthonormal basis. This is well defined because the sum converges absolutely. Indeed,

$$\sum_k |(Se_k, Te_k)| \leq \sum_k |Se_k| |Te_k| \leq \left(\sum_k |Se_k|^2 \right)^{1/2} \left(\sum_k |Te_k|^2 \right)^{1/2} < \infty$$

Definition 22.5.17 For $X \in G$ and $Y \in H$, $X \otimes Y(h) \equiv X(h, Y)$. This is a continuous linear map from H to G because $\|X \otimes Y(h)\| = \|X(h, Y)\| \leq \|h\| \|Y\| \|X\|$. Next we show $X \otimes Y \in \mathcal{L}_2(H, G)$ among other things.

Theorem 22.5.18 $\mathcal{L}_2(H, G)$ is a separable Hilbert space with norm given by

$$\|T\|_{\mathcal{L}_2} \equiv \left(\sum_k \|Te_k\|^2 \right)^{1/2}$$

where $\{e_k\}$ is some fixed orthonormal basis for H . Also $\mathcal{L}_2(H, G) \subseteq \mathcal{L}(H, G)$ and

$$\|T\| \leq \|T\|_{\mathcal{L}_2}. \quad (22.20)$$

All Hilbert Schmidt operators are compact. Also for $X \in G$ and $Y \in H$, $X \otimes Y \in \mathcal{L}_2(H, G)$ and

$$\|X \otimes Y\|_{\mathcal{L}_2} = \|X\|_G \|Y\|_H \quad (22.21)$$

If T is Hilbert Schmidt, then so is T^*T and T^* and ST for any $S \in \mathcal{L}(G, H)$. If $T = T^*$ and $G = H$, then the choice of orthonormal basis in computing $\|T\|_{\mathcal{L}_2}$ is not important.

Proof: Is $\|T\|_{\mathcal{L}_2}$ really a norm? This obviously is so except for the triangle inequality. But this follows from the triangle inequality.

$$\begin{aligned} \|T + S\|_{\mathcal{L}_2} &\equiv \left(\sum_k \|Te_k + Se_k\|^2 \right)^{1/2} \leq \left(\sum_k \|Te_k\|^2 \right)^{1/2} + \left(\sum_k \|Se_k\|^2 \right)^{1/2} \\ &= \|T\|_{\mathcal{L}_2} + \|S\|_{\mathcal{L}_2} \end{aligned}$$

Next is the claim that \mathcal{L}_2 is a Hilbert space. So pick an orthonormal basis $\{e_k\}$. It is clear that \mathcal{L}_2 is an inner product space with respect to the inner product described above in Definition 22.5.16.

Consider completeness. Suppose that $\{T_n\}$ is a Cauchy sequence in $\mathcal{L}_2(H, G)$. Then from 22.20 $\{T_n\}$ is a Cauchy sequence in $\mathcal{L}(H, G)$ and so there exists a unique T such

that $\lim_{n \rightarrow \infty} \|T_n - T\| = 0$. Then it only remains to verify $T \in \mathcal{L}_2(H, G)$. But by Fatou's lemma, for $\{e_k\}$ orthonormal,

$$\sum_k \|Te_k\|^2 \leq \liminf_{n \rightarrow \infty} \sum_k \|T_n e_k\|^2 \equiv \liminf_{n \rightarrow \infty} \|T_n\|_{\mathcal{L}_2}^2 < \infty.$$

All that remains is to verify $\mathcal{L}_2(H, G)$ is separable and these Hilbert Schmidt operators are compact. I will show an orthonormal basis for $\mathcal{L}_2(H, G)$ is $\{f_j \otimes e_k\}$ where $\{f_k\}$ is an orthonormal basis for G and $\{e_k\}$ is an orthonormal basis for H . Here, for $f \in G$ and $e \in H$, $f \otimes e(x) \equiv (x, e)f$.

I need to show $f_j \otimes e_k \in \mathcal{L}_2(H, G)$ and that it is an orthonormal basis for $\mathcal{L}_2(H, G)$ as claimed. Let the $\{e_k\}$ be the orthonormal basis used to define the inner product but the $\{f_j\}$ are just an arbitrary orthonormal basis for G .

$$\sum_k \|f_j \otimes e_i(e_k)\|^2 = \sum_k \|f_j \delta_{ik}\|^2 = \|f_j\|^2 = 1 < \infty$$

so each of these operators is in $\mathcal{L}_2(H, G)$. As noted above, they are also each continuous. Next I show they are orthonormal. From the definition of the inner product,

$$\begin{aligned} (f_j \otimes e_k, f_s \otimes e_r) &= \sum_p (f_j \otimes e_k(e_p), f_s \otimes e_r(e_p)) \\ &= \sum_p \delta_{rp} \delta_{kp} (f_j, f_s) = \sum_p \delta_{rp} \delta_{kp} \delta_{js} \end{aligned}$$

If $j = s$ and $k = r$ this reduces to 1. Otherwise, this gives 0. Thus these operators are orthonormal.

Why is $\mathcal{L}_2(H, G)$ a separable Hilbert space? Let $T \in \mathcal{L}_2(H, G)$. Consider

$$T_n \equiv \sum_{i=1}^n \sum_{j=1}^n (Te_i, f_j) f_j \otimes e_i$$

Then

$$T_n e_k = \sum_{i=1}^n \sum_{j=1}^n (Te_i, f_j) (e_k, e_i) f_j = \sum_{j=1}^n (Te_k, f_j) f_j,$$

a partial sum for Te_k . It follows $\|T_n e_k\| \leq \|Te_k\|$ and $\lim_{n \rightarrow \infty} T_n e_k = Te_k$. Therefore, from the dominated convergence theorem,

$$\lim_{n \rightarrow \infty} \|T - T_n\|_{\mathcal{L}_2}^2 \equiv \lim_{n \rightarrow \infty} \sum_k \|(T - T_n)e_k\|^2 = 0.$$

Therefore, the linear combinations of the $f_j \otimes e_i$ are dense in $\mathcal{L}_2(H, G)$ and this proves completeness of the orthonormal basis.

By only using rational scalars in the linear combinations we see that $\mathcal{L}_2(H, G)$ is separable. From 22.20 it also shows that every $T \in \mathcal{L}_2(H, G)$ is the limit in the operator norm of a sequence of compact operators. This follows because each of the $f_j \otimes e_i$ is easily seen to be a compact operator because if $B \subseteq H$ is bounded, then $(f_j \otimes e_i)(B)$ is a bounded subset of a one dimensional vector space so it is pre-compact. Thus T_n is compact, being a finite sum of these. By Lemma 22.5.10, so is T .

Consider 22.21.

$$\begin{aligned}\|X \otimes Y\|_{\mathcal{L}_2}^2 &\equiv \sum_k \|X \otimes Y(e_k)\|_G^2 \equiv \sum_k \|X(e_k, Y)\|_G^2 \\ &= \|X\|_G^2 \sum_k |(e_k, Y)|^2 = \|X\|_G^2 \|Y\|_H^2\end{aligned}$$

Finally, consider the last claim. I need to show that for self adjoint operators in \mathcal{L}_2 the choice of orthonormal basis does not matter. This is because if $\{e_k\}, \{f_j\}$ are two orthonormal bases, then

$$\sum_k \|Te_k\|^2 = \sum_k \sum_j |(Te_k, f_j)|^2 = \sum_j \sum_k |(e_k, Tf_j)|^2 = \sum_j \|Tf_j\|^2. \blacksquare$$

In fact the orthonormal basis does not matter in defining the norm of any Hilbert Schmidt operator which is not surprising from linear algebra. I will show this as an application a little later in Proposition 22.6.4.

22.6 Roots of Positive Linear Maps

In this section, H will be a Hilbert space, real or complex, and T will denote an operator which satisfies the following definition. This will be a more general result than the above because it will hold for infinite dimensional spaces.

Definition 22.6.1 Let T satisfy $T = T^*$ (Hermitian) and for all $x \in H$,

$$(Tx, x) \geq 0 \quad (22.22)$$

Such an operator is referred to as positive and self adjoint. It is probably better to refer to such an operator as “nonnegative” since the possibility that $Tx = 0$ for some $x \neq 0$ is not being excluded. Instead of “self adjoint” you can also use the term, Hermitian. To save on notation, write $T \geq 0$ to mean T is positive, satisfying 22.22. When we say $A \leq B$ this means $B - A \geq 0$.

A useful theorem about the existence of roots of positive self adjoint operators is presented. This proof is very elementary. I found it in [34] for square roots.

22.6.1 The Product of Positive Self Adjoint Operators

With the above definition here is a fundamental result about positive self adjoint operators.

Proposition 22.6.2 Let S, T be positive and self adjoint such that $ST = TS$. Then ST is also positive and self adjoint.

Proof: It is obvious that ST is self adjoint.

$$(STx, y) = (TSx, y) = (Sx, Ty) = (x, STy)$$

The only problem is to show that ST is positive. The idea is to write $S = S_{n+1} + \sum_{k=0}^n S_k^2$ where $S_0 = S$ and the operators S_k are self adjoint. This is because if you have (TS^2x, x) ,

where everything commutes, this equals $(STSx, x) = (TSx, Sx) \geq 0$. Thus it will be possible to deal with the terms of the sum which are squared. First assume $(Sx, x) \leq (x, x)$ so $S \leq I$.

Define a sequence recursively as follows.

$$S_{n+1} = S_n - S_n^2, \quad S \equiv S_0 \quad (22.23)$$

Then $\sum_{k=0}^n S_k^2 = \sum_{k=0}^n (S_k - S_{k+1}) = S - S_{n+1}$, $S = S_{n+1} + \sum_{k=0}^n S_k^2$. Now $S_0 \geq 0$ by assumption. Assume $S_n \geq 0$. Then

$$S_{n+1} = S_n - S_n^2 = (I - S_n)S_n(S_n + (I - S_n)) = S_n^2(I - S_n) + (I - S_n)^2 S_n$$

It follows that $S_{n+1} \geq 0$ because clearly those two terms on the end are positive. Therefore,

$$(Sx, x) = (S_{n+1}x, x) + \sum_{k=0}^n (S_k^2 x, x) \geq \sum_{k=0}^n \|S_k x\|^2, \quad (Sx, x) \geq \sum_{k=0}^{\infty} \|S_k x\|^2$$

also and so $\lim_{k \rightarrow \infty} \|S_k x\| = 0$. $TSx = TS_{n+1}x + \sum_{k=0}^n TS_k^2 x$.

$$(TSx, x) = (S_{n+1}x, Tx) + \sum_{k=0}^n (TS_k^2 x, x) = (S_{n+1}x, Tx) + \sum_{k=0}^n (TS_k x, S_k x)$$

so passing to a limit as $n \rightarrow \infty$, $(TSx, x) = 0 + \limsup_{n \rightarrow \infty} \sum_{k=0}^n (TS_k x, S_k x) \geq 0$.

Thus if $S \leq I$, the theorem is proved. If S is general, $\frac{S}{\|S\|} \leq I$. In this case, $\left(T \frac{S}{\|S\|} x, x\right) = \left(\frac{S}{\|S\|} Tx, x\right) \geq 0$ and so $(STx, x) \geq 0$. ■

The proposition is like the familiar statement about real numbers which says that when you multiply two nonnegative real numbers the result is a nonnegative real number.

22.6.2 Roots of Positive Self Adjoint Operators

With this preparation, it is time to give the theorem about roots.

Theorem 22.6.3 *Let $T \in \mathcal{L}(H, H)$ be a positive self adjoint linear operator. Then for $m \in \mathbb{N}$, there exists a unique m^{th} root A with the following properties. $A^m = T$, A is positive and self adjoint, A commutes with every operator which commutes with T .*

Proof: Define the following sequence of operators:

$$A_0 \equiv 0, \quad A_{n+1} \equiv A_n + \frac{1}{m} (T - A_n^m)$$

Say $T \leq I$.

Claim 1: $A_n \leq I$.

Proof of Claim 1: True if $n = 0$. Assume true for n . Then

$$\begin{aligned} I - A_{n+1} &= I - A_n + \frac{1}{m} (A_n^m - T) \geq I - A_n + \frac{1}{m} (A_n^m - I) \\ &= I - A_n - \frac{1}{m} (I - A_n^m) \\ &= (I - A_n) - \frac{1}{m} (I - A_n) (I + \cdots + A_n^{m-1}) \end{aligned}$$

Now, since $A_n \leq I, I + \cdots + A_n^{m-1} \leq mI$, it follows that

$$= (I - A_m) \left(I - \frac{1}{m} (I + \cdots + A_n^{m-1}) \right) \geq (I - A_m) (I - I) = 0$$

so by induction, $A_n \leq I$.

Claim 2: $A_n \leq A_{n+1}$.

Proof of Claim 2: From the definition of A_n , this is true if $n = 0$ because

$$A_1 = T \geq 0 = A_0.$$

Suppose true for n . Then from Claim 1,

$$\begin{aligned} A_{n+2} - A_{n+1} &= A_{n+1} + \frac{1}{m} (T - A_{n+1}^m) - \left[A_n + \frac{1}{m} (T - A_n^m) \right] \\ &= A_{n+1} - A_n + \frac{1}{m} (A_n^m - A_{n+1}^m) \\ &= (A_{n+1} - A_n) - (A_{n+1} - A_n) \frac{1}{m} (A_{n+1}^{m-1} + A_{n+1}^{m-2} A_n + \cdots + A_n^{m-1}) \\ &\geq (A_{n+1} - A_n) - (A_{n+1} - A_n) I = 0 \end{aligned}$$

since each $A_n, A_{n+1} \leq I$, so this proves the claim.

Claim 3: $A_n \geq 0$

Proof of Claim 3: This is true if $n = 0$. Suppose it is true for n .

$$\begin{aligned} (A_{n+1}x, x) &= (A_nx, x) + \frac{1}{m} (Tx, x) - \frac{1}{m} (A_n^m x, x) \\ &\geq (A_nx, x) + \frac{1}{m} (Tx, x) - \frac{1}{m} (A_nx, x) \geq 0 \end{aligned}$$

because by Proposition 22.6.2, $A_n - A_n^m = A_n (I - A_n^{m-1}) \geq 0$ because $A_n \leq I$.

Thus (A_nx, x) is increasing and bounded above so it converges. Now let $n > k$. Using Proposition 22.6.2 $A_n A_k \geq A_k^2$ and also

$$(A_n - A_k)(A_n + A_k) \leq 2(A_n - A_k).$$

Thus the following holds.

$$\begin{aligned} \|A_nx - A_kx\|^2 &= ((A_n - A_k)^2 x, x) = (A_n^2 x, x) - 2(A_n A_k x, x) + (A_k^2 x, x) \\ &\leq (A_n^2 x, x) - 2(A_k^2 x, x) + (A_k^2 x, x) = ((A_n - A_k)(A_n + A_k)x, x) \\ &\leq 2[(A_nx, x) - (A_kx, x)] \end{aligned}$$

which converges to 0 as $k, n \rightarrow \infty$. Therefore, $\lim_{n \rightarrow \infty} A_n x$ exists since $\{A_n x\}$ is a Cauchy sequence. Let this limit be Ax . Then clearly A is linear. Also, since each $A_n \geq 0$ and self adjoint, the Cauchy Schwarz inequality implies

$$|(Ax, y)| = \lim_{n \rightarrow \infty} |(A_n x, y)| \leq \limsup_{n \rightarrow \infty} \left| (A_n x, x)^{1/2} (A_n y, y)^{1/2} \right| \leq \|x\| \|y\|$$

so A is also continuous. Now $(Ax, x) = \lim_{n \rightarrow \infty} (A_n x, x) \geq 0$ so A is positive and it is clearly also self adjoint since each A_n is. From passing to the limit in the definition of A_n ,

$$Ax = Ax + \frac{1}{m} (Tx - A^m x)$$

and so $Tx = A^m x$. This proves the theorem in the case that $T \leq I$. Then if $T > I$, consider $T/\|T\|$. $T/\|T\| \leq I$ and so there is B such that $B^m = T/\|T\|$. Let $A = \|T\|^{1/m} B$. This proves the existence of the m^{th} root. It is clear that A commutes with every continuous linear operator that commutes with T because this is true of each of the iterates. In fact, each of these is just a polynomial in T . It remains to verify uniqueness.

Next suppose both A and B are m^{th} roots of T having all the properties stated in the theorem. Then $AB = BA$ because both A and B commute with every operator which commutes with T . Then from Proposition 22.6.2,

$$((A^{m-1} + A^{m-2}B + \dots + B^{m-1})(A - B)x, (A - B)x) \geq 0 \quad (22.24)$$

Therefore, $((A^m - B^m)x, (A - B)x) = (0, (A - B)x) = 0$.

Now this means $(A^k B^l (A - B)x, (A - B)x) = 0$ for all $k + l = m - 1$ since the sum of such terms is 0 and each of them is nonnegative. Now this implies

$$(\sqrt{A^k B^l} (A - B)x, \sqrt{A^k B^l} (A - B)x) = 0$$

and so $\sqrt{A^k B^l} (A - B)x = 0 \Rightarrow A^k B^l (A - B)x = 0, k + l = m - 1$. Then, using the binomial theorem,

$$0 = \sum_{j=0}^{m-1} \binom{m-1}{j} A^{m-1-j} B^j (-1)^j (A - B)x = (A - B)^m x$$

This clearly implies $A = B$. To see this, consider $m = 7$.

If $m = 7$, $(A - B)^7 x = 0$ so

$$(A - B)^8 x = 0$$

so $((A - B)^4 x, (A - B)^4 x) = 0$ which implies $(A - B)^4 x = 0$ which implies

$$((A - B)^2 x, (A - B)^2 x) = 0$$

so $((A - B)^2 x, x) = 0$ which yields $((A - B)x, (A - B)x) = 0$, so $(A - B) = 0$. ■

This next was shown earlier, but this is a nice way to think of it in terms of a square root.

Proposition 22.6.4 Let $T \in \mathcal{L}_2(H, G)$. Then if $\{e_k\}, \{f_j\}$ are two orthonormal bases, then $\sum_k \|Te_k\|^2 = \sum_k \|Tf_k\|^2$.

Proof: T^*T is self adjoint and in $\mathcal{L}_2(H, H)$. Therefore,

$$\sum_k \left\| \sqrt{T^*T} e_k \right\|^2 = \sum_k (T^*T e_k, e_k) = \sum_k \|Te_k\|^2$$

which is finite. Thus $\sqrt{T^*T}$ is self adjoint and in $\mathcal{L}_2(H, H)$ and so

$$\sum_k \left\| \sqrt{T^*T} e_k \right\|^2 = \sum_k \left\| \sqrt{T^*T} f_k \right\|^2 = \sum_k \|T f_k\|^2$$

showing that $\sum_k \|T e_k\|^2 = \sum_k \|T f_k\|^2$. ■

There is a whole book on powers of operators. This has just given a short introduction. See [33].

22.7 Differential Equations in Banach Space

Here we consider the initial value problem for functions which have values in a Banach space. Let X be a Banach space.

Definition 22.7.1 Define $BC([a, b]; X)$ as bounded continuous functions f which have values in the Banach space X . For $f \in BC([a, b]; X)$, γ a real number. Then

$$\|f\|_\gamma \equiv \sup_{t \in [a, b]} \left\| f(t) e^{\gamma(t-a)} \right\| \quad (22.25)$$

Then this is a norm. The usual norm is given by $\|f\| \equiv \sup_{t \in [a, b]} \|f(t)\|$.

Lemma 22.7.2 $\|\cdot\|_\gamma$ is a norm for $BC([a, b]; X)$ and $BC([a, b]; X)$ is a complete normed linear space. Also, a sequence is Cauchy in $\|\cdot\|_\gamma$ if and only if it is Cauchy in $\|\cdot\|$.

Proof: First consider the claim about $\|\cdot\|_\gamma$ being a norm. To simplify notation, let $T = [a, b]$. It is clear that $\|f\|_\gamma = 0$ if and only if $f = 0$ and $\|f\|_\gamma \geq 0$. Also,

$$\|\alpha f\|_\gamma \equiv \sup_{t \in T} \left\| \alpha f(t) e^{\gamma(t-a)} \right\| = |\alpha| \sup_{t \in T} \left\| f(t) e^{\gamma(t-a)} \right\| = |\alpha| \|f\|_\gamma$$

so it does what is should for scalar multiplication. Next consider the triangle inequality.

$$\begin{aligned} \|f + g\|_\gamma &= \sup_{t \in T} \left\| (f(t) + g(t)) e^{\gamma(t-a)} \right\| \leq \sup_{t \in T} \left(\left\| f(t) e^{\gamma(t-a)} \right\| + \left\| g(t) e^{\gamma(t-a)} \right\| \right) \\ &\leq \sup_{t \in T} \left\| f(t) e^{\gamma(t-a)} \right\| + \sup_{t \in T} \left\| g(t) e^{\gamma(t-a)} \right\| = \|f\|_\gamma + \|g\|_\gamma \end{aligned}$$

The rest follows from the next inequalities.

$$\begin{aligned} \|f\| &\equiv \sup_{t \in T} \|f(t)\| = \sup_{t \in T} \left\| f(t) e^{\gamma(t-a)} e^{-\gamma(t-a)} \right\| \leq e^{|\gamma(b-a)|} \|f\|_\gamma \\ &\equiv e^{|\gamma(b-a)|} \sup_{t \in T} \left\| f(t) e^{\gamma(t-a)} \right\| \leq \left(e^{|\gamma(b-a)|} \right)^2 \sup_{t \in T} \|f(t)\| = \left(e^{|\gamma(b-a)|} \right)^2 \|f\| \quad \blacksquare \end{aligned}$$

Now consider the ordinary initial value problem

$$x'(t) = F(t, x(t)), \quad x(t_0) = x_0, \quad t \in [a, b], \quad t_0 \in [a, b] \quad (22.26)$$

where here $F : [a, b] \times X \rightarrow X$ is continuous and satisfies the Lipschitz condition

$$\|F(t, x) - F(t, y)\| \leq K \|x - y\|, \quad F : [a, b] \times X \rightarrow X \text{ is continuous} \quad (22.27)$$

Thanks to the fundamental theorem of calculus, there exists a solution to 22.26 if and only if it is a solution to the integral equation

$$x(t) = x_0 + \int_{t_0}^t F(s, x(s)) ds \quad (22.28)$$

Then we have the following theorem.

Theorem 22.7.3 *Let 22.27 hold. Then there exists a unique solution to 22.26 in $BC([a, b]; X)$.*

Proof: Use the norm of 22.25 where $\gamma \neq 0$ is described later. Let $T : BC([a, b]; X) \rightarrow BC([a, b]; X)$ be defined by

$$Tx(t) \equiv x_0 + \int_{t_0}^t F(s, x(s)) ds$$

Then

$$\begin{aligned} \|Tx(t) - Ty(t)\|_X &= \left\| \int_{t_0}^t F(s, x(s)) ds - \int_{t_0}^t F(s, y(s)) ds \right\| \\ &\leq K \int_{t_0}^t \|x(s) - y(s)\| ds = K \int_{t_0}^t \left\| (x(s) - y(s)) e^{\gamma(s-a)} e^{-\gamma(s-a)} \right\| ds \\ &\leq K \int_{t_0}^t e^{-\gamma(s-a)} ds \|x - y\|_\gamma = K \left(\frac{e^{-\gamma(t-a)}}{-\gamma} + \frac{e^{-\gamma(t_0-a)}}{\gamma} \right) \|x - y\|_\gamma \end{aligned}$$

Therefore, letting $\gamma < 0$

$$e^{\gamma(t-a)} \|Tx(t) - Ty(t)\|_X \leq K \left(\frac{1}{-\gamma} + \frac{e^{\gamma(t-t_0)}}{\gamma} \right) \|x - y\|_\gamma < K \left(\frac{1}{|\gamma|} \right) \|x - y\|_\gamma$$

$$\|Tx - Ty\|_\gamma \leq K \left(\frac{1}{|\gamma|} \right) \|x - y\|_\gamma$$

Letting $\gamma = -m^2$, this reduces to

$$\|Tx - Ty\|_{-m^2} \leq \frac{K}{m^2} \|x - y\|_{-m^2}$$

and so if $K/m^2 < 1/2$, this shows the solution to the integral equation is the unique fixed point of a contraction mapping defined on $BC([a, b]; X)$. This shows existence and uniqueness of the initial value problem 22.26. ■

Definition 22.7.4 *Let $S : [0, \infty) \rightarrow \mathcal{L}(X, X)$ be continuous and satisfy*

1. $S(t+s) = S(t)S(s)$ called the semigroup identity.
2. $S(0) = I$
3. $\lim_{h \rightarrow 0+} \frac{S(h)x - x}{h} = Ax$ for A a densely defined closed linear operator whenever $x \in D(A) \subseteq X$.

Then S is called a continuous semigroup and A is said to generate S .

Then we have the following corollary of Theorem 22.7.3. First note the following. For $t \geq 0$ and $h \geq 0$, if $x \in D(A)$, the semigroup identity implies

$$\lim_{h \rightarrow 0} \frac{S(t+h)x - S(t)x}{h} = \lim_{h \rightarrow 0} S(t) \frac{S(h)x - x}{h} = S(t) \lim_{h \rightarrow 0} \frac{S(h)x - x}{h} \equiv S(t)Ax$$

As shown above, $\mathcal{L}(X, X)$ is a Banach space with the operator norm whenever X is a Banach space.

Corollary 22.7.5 *Let X be a Banach space and let $A \in \mathcal{L}(X, X)$. Let $S(t)$ be the solution in $\mathcal{L}(X, X)$ to*

$$S'(t) = AS(t), \quad S(0) = I, \quad t \geq 0 \quad (22.29)$$

Then $t \rightarrow S(t)$ is a continuous semigroup whose generator is A . In this case that A is actually defined on all of X , not just on a dense subset. Furthermore, in this case where $A \in \mathcal{L}(X, X)$, $S(t)A = AS(t)$. If $T(t)$ is any semigroup having A as a generator, then $T(t) = S(t)$. Also you can express $S(t)$ as a power series, $S(t) = \sum_{n=0}^{\infty} \frac{(At)^n}{n!}$.

Proof: The solution to the initial value problem 22.29 exists on $[-b, b]$ for all b so it exists on all of \mathbb{R} thanks to the uniqueness on every finite interval. First consider the semigroup property. Let $\Psi(t) \equiv S(t+s)$, $\Phi(t) \equiv S(t)S(s)$. Then

$$\Psi'(t) = S'(t+s) = AS(t+s) = A\Psi(t), \quad \Psi(0) = S(s)$$

$$\Phi'(t) = S'(t)S(s) = AS(t)S(s) = A\Phi(t), \quad \Phi(0) = S(s)$$

By uniqueness, $\Phi(t) = \Psi(t)$ for all $t \geq 0$. Thus $S(t)S(s) = S(t+s) = S(s)S(t)$. Now from this, for $t > 0$

$$S(t)A = S(t) \lim_{h \rightarrow 0} \frac{S(h) - I}{h} = \lim_{h \rightarrow 0} S(t) \frac{S(h) - I}{h} = \lim_{h \rightarrow 0} \frac{S(h) - I}{h} S(t) = AS(t).$$

As to A being the generator of $S(t)$, letting $x \in X$, then from the differential equation solved,

$$\lim_{h \rightarrow 0+} \frac{S(h)x - x}{h} = \lim_{h \rightarrow 0+} \frac{1}{h} \int_0^h AS(t)x dt = AS(0)x = Ax.$$

If $T(t)$ is a semigroup generated by A then for $t > 0$,

$$T'(t) \equiv \lim_{h \rightarrow 0} \frac{T(t+h) - T(t)}{h} = \lim_{h \rightarrow 0} \frac{T(h) - I}{h} T(t) = AT(t)$$

and $T(0) = I$. However, uniqueness applies because T and S both satisfy the same initial value problem and this yields $T(t) = S(t)$.

To show the power series equals $S(t)$ it suffices to show it satisfies the initial value problem. Using the mean value theorem,

$$\sum_{n=0}^{\infty} \frac{A^n((t+h)^n - t^n)}{n!} = \sum_{n=1}^{\infty} \frac{A^n(t + \theta_n(h))^{n-1}}{(n-1)!}$$

where $\theta_n(h) \in (0, h)$. Then taking a limit as $h \rightarrow 0$ and using the dominated convergence theorem, the limit of the difference quotient is

$$\sum_{n=1}^{\infty} \frac{A^n t^{n-1}}{(n-1)!} = A \sum_{n=1}^{\infty} \frac{A^{n-1} t^{n-1}}{(n-1)!} = A \sum_{n=0}^{\infty} \frac{(At)^n}{n!}$$

Thus $\sum_{n=0}^{\infty} \frac{(At)^n}{n!}$ satisfies the differential equation. It clearly satisfies the initial condition. Hence it equals $S(t)$. ■

Note that as a consequence of the above argument showing that T and S are the same, it follows that $T(t)A = AT(t)$ so one obtains that if the generator is a bounded linear operator, then the semigroup commutes with this operator.

When dealing with differential equations, one of the best tools is Gronwall's inequality. This is presented next.

Theorem 22.7.6 Suppose u is nonnegative, continuous, and real valued and that

$$u(t) \leq C + \int_0^t ku(s) ds, \quad k \geq 0$$

Then $u(t) \leq Ce^{kt}$.

Proof: Let $w(t) \equiv \int_0^t ku(s) ds$. Then

$$w'(t) = ku(t) \leq kC + kw(t)$$

and so $w'(t) - kw(t) \leq kC$ which implies $\frac{d}{dt}(e^{-kt}w(t)) \leq kCe^{-kt}$. Therefore,

$$e^{-kt}w(t) \leq Ck \int_0^t e^{-ks} ds = Ck \left(\frac{1}{k} - \frac{1}{k}e^{-kt} \right)$$

so $w(t) \leq C(e^{kt} - 1)$. From the original inequality, $u(t) \leq C + w(t) \leq C + C(e^{kt} - 1) = Ce^{kt}$. ■

22.8 General Theory of Continuous Semigroups

Much more on semigroups is available in Yosida [60]. This is just an introduction to the subject.

22.8.1 Generators of Semigroups

Definition 22.8.1 A strongly continuous semigroup defined on X , a Banach space is a function $S: [0, \infty) \rightarrow \mathcal{L}(X, X)$ which satisfies the following for all $x_0 \in X$.

$$\begin{aligned} S(t) &\in \mathcal{L}(X, X), S(t+s) = S(t)S(s), \\ t &\rightarrow S(t)x_0 \text{ is continuous, } \lim_{t \rightarrow 0+} S(t)x_0 = x_0 \end{aligned}$$

Sometimes such a semigroup is said to be C_0 . It is said to have the linear operator A as its generator if

$$D(A) \equiv \left\{ x : \lim_{h \rightarrow 0} \frac{S(h)x - x}{h} \text{ exists} \right\}$$

and for $x \in D(A)$, A is defined by

$$\lim_{h \rightarrow 0} \frac{S(h)x - x}{h} \equiv Ax$$

The assertion that $t \rightarrow S(t)x_0$ is continuous and that $S(t) \in \mathcal{L}(X, X)$ is not sufficient to say there is a bound on $\|S(t)\|$ for all $t \geq 0$. Also the assertion that for each x_0 , $\lim_{t \rightarrow 0+} S(t)x_0 = x_0$ is not the same as saying that $S(t) \rightarrow I$ in $\mathcal{L}(X, X)$. It is a much weaker assertion. The next theorem gives information on the growth of $\|S(t)\|$. It turns out it has exponential growth. Thus $S(t)$ is a lot like e^t .

Lemma 22.8.2 *Let $M \equiv \sup \{\|S(t)\| : t \in [0, T]\}$. Then $M < \infty$.*

Proof: If this is not true, then there exists $t_n \in [0, T]$ such that $\|S(t_n)\| \geq n$. That is the operators $S(t_n)$ are not uniformly bounded. By the uniform boundedness principle, Theorem 21.1.9, there exists $x \in X$ such that $\|S(t_n)x\|$ is not bounded. However, this is impossible because it is given that $t \rightarrow S(t)x$ is continuous on $[0, T]$ and so $t \rightarrow \|S(t)x\|$ must achieve its maximum on this compact set. ■

Now here is the main result for growth of $\|S(t)\|$.

Theorem 22.8.3 *For M described in Lemma 22.8.2, there exists α such that*

$$\|S(t)\| \leq Me^{\alpha t}, t \geq 0$$

In fact, α can be chosen such that $M^{1/T} = e^\alpha$.

Proof: Let t be arbitrary. Then $t = mT + r(t)$ where $0 \leq r(t) < T$. Then by the semigroup property

$$\|S(t)\| = \|S(mT + r(t))\| = \|S(r(t))S(T)^m\| \leq M^{m+1}$$

Now $mT \leq t \leq mT + r(t) \leq (m+1)T$ and so $m \leq \frac{t}{T} \leq m+1$. Therefore,

$$\|S(t)\| \leq M^{(t/T)+1} = M \left(M^{1/T}\right)^t.$$

Let $M^{1/T} \equiv e^\alpha$ and then $\|S(t)\| \leq Me^{\alpha t}$ ■

Definition 22.8.4 *Let $S(t)$ be a continuous semigroup as described above. It is called a contraction semigroup if for all $t \geq 0$, $\|S(t)\| \leq 1$. It is called a bounded semigroup if there exists M such that for all $t \geq 0$, $\|S(t)\| \leq M$.*

Note that for $S(t)$ an arbitrary continuous semigroup satisfying $\|S(t)\| \leq Me^{\alpha t}$, It follows that the semigroup, $T(t) = e^{-\alpha t}S(t)$ is a bounded semigroup which satisfies $\|T(t)\| \leq M$.

The next proposition has to do with taking a Laplace transform of a semigroup. It suffices to let the integral be the usual Riemann integral for a function having values in a Banach space. You define it the same way as in single variable calculus in terms of limits of Riemann sums. Later, this will be generalized.

Proposition 22.8.5 *Given a continuous semigroup $S(t)$, its generator A exists and is a closed densely defined operator. Furthermore, for $\|S(t)\| \leq Me^{\alpha t}$ and $\lambda > \alpha$, $\lambda I - A$ is one to one and onto from $D(A)$ to X . Also $(\lambda I - A)^{-1}$ maps X onto $D(A)$ and is in $\mathcal{L}(X, X)$. Also for these values of $\lambda > \alpha$, $(\lambda I - A)^{-1}x = \int_0^\infty e^{-\lambda t} S(t)x dt$. For $\lambda > \alpha$, the following estimate holds.*

$$\|(\lambda I - A)^{-1}\| \leq \frac{M}{|\lambda - \alpha|} \quad (22.30)$$

Proof: First note $D(A) \neq \emptyset$. In fact $0 \in D(A)$. It follows from Theorem 22.8.3 that for all λ larger than α , one can define a Laplace transform, $R(\lambda)x \equiv \int_0^\infty e^{-\lambda t} S(t)x dt \in X$. The integral is the ordinary improper Riemann integral. Note that for $\lambda > \alpha$, $R(\lambda) \in \mathcal{L}(X, X)$ thanks to the estimates. Indeed, approximating with Riemann sums, to justify the details,

$$\|R(\lambda)x\| \leq \int_0^\infty e^{-\lambda t} \|S(t)x\| dt \leq \int_0^\infty Me^{-(\lambda-\alpha)t} dt \|x\| \leq \frac{M}{|\lambda - \alpha|} \|x\| \quad (22.31)$$

Claim 1: For $\lambda > \alpha$, $R(\lambda)x \in D(A)$ and $x = (\lambda I - A)R(\lambda)x$ so $R(\lambda)$ is a right inverse of $(\lambda I - A)$.

Proof of Claim 1: From the semigroup formula,

$$\begin{aligned} \frac{S(h)R(\lambda)x - R(\lambda)x}{h} &= \frac{e^{h\lambda} \int_0^\infty e^{-\lambda(t+h)} S(t+h)x dt - \int_0^\infty e^{-\lambda t} S(t)x dt}{h} = \\ &= \frac{e^{h\lambda} \int_h^\infty e^{-\lambda t} S(t)x dt - \int_0^\infty e^{-\lambda t} S(t)x dt}{h} = \frac{(e^{h\lambda} - 1)R(\lambda)x - e^{\lambda h} \int_0^h e^{-\lambda t} S(t)x dt}{h} \end{aligned}$$

Then it follows that the limit as $h \rightarrow 0$ exists and equals $\lambda R(\lambda)x - x$ which by definition of A is $A(R(\lambda)x)$. So by definition, $R(\lambda)x \in D(A)$ as claimed, and $\lambda IR(\lambda)x - A(R(\lambda)x) = x$ and so $x = (\lambda I - A)R(\lambda)x$. This shows **Claim 1**.

Claim 2: $D(A)$ is dense in X and for any $x \in X$, $\lim_{\lambda \rightarrow \infty} \lambda R(\lambda)x = x$.

Proof of Claim 2: Note that $\int_0^\infty \lambda e^{-\lambda t} dt = 1$ and so

$$\begin{aligned} \|\lambda R(\lambda)x - x\| &= \left\| \int_0^\infty \lambda e^{-\lambda t} S(t)x dt - x \right\| = \left\| \int_0^\infty \lambda e^{-\lambda t} (S(t)x - x) dt \right\| \\ &\leq \int_0^\infty \lambda e^{-\lambda t} \|S(t)x - x\| dt \end{aligned}$$

which from the estimates and standard approximate identity type arguments converges to 0 as follows: Let $\varepsilon > 0$ be given. Then choose δ such that $\|S(t)x - x\| < \varepsilon$ if $0 \leq t \leq \delta$. Then for λ large enough the second term in the following is no more than ε

$$\int_0^\infty \lambda e^{-\lambda t} \|S(t)x - x\| dt \leq \int_0^\delta \lambda e^{-\lambda t} \varepsilon dt + \int_\delta^\infty \lambda e^{-\lambda t} (Me^{\alpha t} + 1) \|x\| dt.$$

Thus for λ large enough, $\int_0^\infty \lambda e^{-\lambda t} \|S(t)x - x\| dt < 2\varepsilon$. This shows that $D(A)$ is dense in X and for any x , $\lim_{\lambda \rightarrow \infty} \lambda R(\lambda)x = x$. This proves **Claim 2**.

Claim 3: For $\lambda > \alpha$, $x = R(\lambda)(\lambda I - A)x$ for $x \in D(A)$ so $(\lambda I - A)$ is one to one and $R(\lambda)$ is a left inverse also. Thus $R(\lambda) = (\lambda I - A)^{-1}$ and from 22.31, estimate 22.30 holds.

Proof of Claim 3: If $x \in D(A)$, you could approximate with Riemann sums and pass to a limit and obtain the following for $\lambda > \alpha$.

$$\left\| R(\lambda) \left(\frac{S(h)x - x}{h} \right) - R(\lambda)Ax \right\| = \left\| \int_0^\infty e^{-\lambda t} S(t) \left(\frac{S(h)x - x}{h} - Ax \right) dt \right\|$$

$$\leq \int_0^\infty \|e^{-\lambda t} S(t)\| \left\| \frac{S(h)x - x}{h} - Ax \right\| dt$$

Then, passing to a limit as $h \rightarrow 0+$, this integrand converges uniformly to 0 so for all $\lambda > \alpha$,

$$\lim_{h \rightarrow 0} R(\lambda) \left(\frac{S(h)x - x}{h} \right) = R(\lambda)Ax \quad (22.32)$$

Also, $S(h)$ commutes with $R(\lambda)$. This follows from approximating with Riemann sums and taking a limit. Thus also

$$\lim_{h \rightarrow 0} R(\lambda) \left(\frac{S(h)x - x}{h} \right) = \lim_{h \rightarrow 0} \left(\frac{S(h)R(\lambda)x - R(\lambda)x}{h} \right) = AR(\lambda)x$$

so we have for $x \in D(A)$, $R(\lambda)Ax = AR(\lambda)x$. However, this implies

$$R(\lambda)(\lambda I - A)x = (\lambda I - A)R(\lambda)x = x$$

from **Claim 1**. Thus $R(\lambda)$ is a left inverse of $(\lambda I - A)$. Since $R(\lambda) = (\lambda I - A)^{-1}$, this shows the estimate 22.30 from 22.31. This proves **Claim 3**.

Why is A a closed operator? Suppose $x_n \rightarrow x$ where $x_n \in D(A)$ and that $Ax_n \rightarrow \xi$. I need to show that this implies that $x \in D(A)$ and that $Ax = \xi$. Thus $x_n \rightarrow x$ and for $\lambda > \alpha$, $(\lambda I - A)x_n \rightarrow \lambda x - \xi$. However, 22.30 shows that $(\lambda I - A)^{-1} = R(\lambda)$ is continuous and so

$$x_n \rightarrow (\lambda I - A)^{-1}(\lambda x - \xi) = x$$

It follows that $x \in D(A)$. Then doing $(\lambda I - A)$ to both sides of the equation, $\lambda x - \xi = \lambda x - Ax$ and so $Ax = \xi$ showing that A is a closed operator as claimed. ■

Definition 22.8.6 The linear mapping for $\lambda > \alpha$ where $\|S(t)\| \leq Me^{\alpha t}$ given by $(\lambda I - A)^{-1} = R(\lambda)$ is called the resolvent.

The following corollary is also very interesting.

Corollary 22.8.7 Let $S(t)$ be a continuous semigroup and let A be its generator. Then for $0 < a < b < \infty$ and $x \in D(A)$, $S(b)x - S(a)x = \int_a^b S(t)Ax dt$ and also for $t > 0$ you can take the derivative from the left,

$$\lim_{h \rightarrow 0+} \frac{S(t)x - S(t-h)x}{h} = S(t)Ax$$

Proof: Letting $y^* \in X'$, you can take y^* inside the integral by approximating with Riemann sums. Thus

$$y^* \left(\int_a^b S(t)Ax dt \right) = \int_a^b y^* \left(S(t) \lim_{h \rightarrow 0} \frac{S(h)x - x}{h} \right) dt$$

The difference quotients are bounded because they converge to Ax . Therefore, from the dominated convergence theorem and using the semigroup property,

$$y^* \left(\int_a^b S(t)Ax dt \right) = \lim_{h \rightarrow 0} \int_a^b y^* \left(S(t) \frac{S(h)x - x}{h} \right) dt$$

$$\begin{aligned}
&= \lim_{h \rightarrow 0} \left(\frac{1}{h} \int_{a+h}^{b+h} y^* S(t) x dt - \frac{1}{h} \int_a^b y^* S(t) x dt \right) \\
&= \lim_{h \rightarrow 0} \left(\frac{1}{h} \int_b^{b+h} y^* S(t) x dt - \frac{1}{h} \int_a^{a+h} y^* S(t) x dt \right) = y^* (S(b)x - S(a)x)
\end{aligned}$$

Since y^* is arbitrary, this proves the first part. Now from what was just shown, if $t > 0$ and h is small enough,

$$\frac{S(t)x - S(t-h)x}{h} = \frac{1}{h} \int_{t-h}^t S(s) A x ds$$

which converges to $S(t)Ax$ as $h \rightarrow 0+$. ■

22.8.2 Hille Yosida Theorem

Given a closed densely defined operator, when is it the generator of a continuous semigroup? This is answered in the following theorem which is called the Hille Yosida theorem. It concerns the case of a bounded semigroup. However, if you have an arbitrary continuous semigroup, $S(t)$, then it was shown above that $S(t)e^{-\alpha t}$ is bounded for suitable α so the case discussed below is obtained.

Theorem 22.8.8 Suppose A is a densely defined linear operator which has the property that for all $\lambda > 0$,

$$(\lambda I - A)^{-1} \in \mathcal{L}(X, X)$$

which means that $\lambda I - A : D(A) \rightarrow X$ is one to one and onto with continuous inverse. Suppose also that for all $n \in \mathbb{N}$,

$$\left\| \left((\lambda I - A)^{-1} \right)^n \right\| \leq \frac{M}{\lambda^n}. \quad (22.33)$$

Then there exists a continuous semigroup $S(t)$ which has A as its generator and satisfies $\|S(t)\| \leq M$ and A is closed. In fact letting

$$S_\lambda(t) \equiv \exp \left(-\lambda + \lambda^2 (\lambda I - A)^{-1} \right) \equiv \exp(A_\lambda)$$

it follows $\lim_{\lambda \rightarrow \infty} S_\lambda(t)x = S(t)x$ uniformly on finite intervals. Conversely, if A is the generator of $S(t)$, a bounded continuous semigroup having $\|S(t)\| \leq M$, then $(\lambda I - A)^{-1} \in \mathcal{L}(X, X)$ for all $\lambda > 0$ and 22.33 holds.

Proof: The condition 22.33 implies, that $\left\| (\lambda I - A)^{-1} \right\| \leq \frac{M}{\lambda}$.

Consider, for $\lambda > 0$, the operator which is defined on $D(A)$, $\lambda(\lambda I - A)^{-1}A$. On $D(A)$, this equals

$$-\lambda I + \lambda^2 (\lambda I - A)^{-1} \quad (22.34)$$

because

$$\begin{aligned}
(\lambda I - A) \lambda (\lambda I - A)^{-1} A &= \lambda A \\
(\lambda I - A) \left(-\lambda I + \lambda^2 (\lambda I - A)^{-1} \right) &= -\lambda (\lambda I - A) + \lambda^2 = \lambda A
\end{aligned}$$

and, by assumption, $(\lambda I - A)$ is one to one. From the second line of 22.34, the operator $-\lambda I + \lambda^2 (\lambda I - A)^{-1}$ makes sense on all of X not just on $D(A)$. Also

$$\left(-\lambda I + \lambda^2 (\lambda I - A)^{-1}\right) (\lambda I - A) = -\lambda (\lambda I - A) + \lambda^2 I = \lambda A$$

$$\lambda A (\lambda I - A)^{-1} (\lambda I - A) = \lambda A$$

so, since $(\lambda I - A)$ is onto, it follows that on X ,

$$-\lambda I + \lambda^2 (\lambda I - A)^{-1} = A \lambda (\lambda I - A)^{-1} \equiv A_\lambda$$

Denote this as A_λ to save notation. Thus on $D(A)$,

$$\lambda A (\lambda I - A)^{-1} = \lambda (\lambda I - A)^{-1} A = A_\lambda$$

although the $\lambda (\lambda I - A)^{-1} A$ only makes sense on $D(A)$. This is summarized next.

Lemma 22.8.9 *There is a bounded linear operator given for $\lambda > 0$ by*

$$-\lambda I + \lambda^2 (\lambda I - A)^{-1} = \lambda A (\lambda I - A)^{-1} \equiv A_\lambda$$

On $D(A)$, $A_\lambda = \lambda (\lambda I - A)^{-1} A$. Also, for all $x \in X$,

$$\lim_{\lambda \rightarrow \infty} \lambda (\lambda I - A)^{-1} x - x = 0. \quad (22.35)$$

Replacing x with Ax , it follows that for all $x \in D(A)$,

$$\lim_{\lambda \rightarrow \infty} A_\lambda x = Ax. \quad (22.36)$$

Proof: First assume $x \in D(A)$

$$\begin{aligned} & \left\| \lambda (\lambda I - A)^{-1} x - x \right\| = \left\| (\lambda I - A)^{-1} (\lambda x - (\lambda I - A)x) \right\| \\ &= \left\| (\lambda I - A)^{-1} Ax \right\| \leq \frac{M}{\lambda} \|Ax\| \end{aligned} \quad (22.37)$$

which converges to 0 as $\lambda \rightarrow \infty$.

Now let x be general and let $\hat{x} \in D(A)$. From 22.33, $\left| \lambda (\lambda I - A)^{-1} \right| \leq M$. Then

$$\begin{aligned} \left\| \lambda (\lambda I - A)^{-1} x - x \right\| &\leq \left\| \lambda (\lambda I - A)^{-1} x - \lambda (\lambda I - A)^{-1} \hat{x} \right\| \\ &\quad + \left\| \lambda (\lambda I - A)^{-1} \hat{x} - \hat{x} \right\| + \|\hat{x} - x\| \end{aligned}$$

Let \hat{x} be close enough to x that the first and last terms on the right added together are less than $\varepsilon/2$. Then whenever λ is large enough, the first part of the argument shows that the middle term is no more than $\varepsilon/2$. This verifies 22.35, 22.36. ■

Now from Corollary 22.7.5, there exists an approximate continuous semigroup $S_\lambda(t)$ generated by A_λ which is the solution to

$$S'_\lambda(t) = A_\lambda S_\lambda(t), S_\lambda(0) = I \quad (22.38)$$

In terms of power series,

$$S_\lambda(t) \equiv e^{-\lambda t} \sum_{k=0}^{\infty} \frac{t^k (\lambda^2 (\lambda I - A)^{-1})^k}{k!} = e^{t(-\lambda I + \lambda^2 (\lambda I - A)^{-1})} \quad (22.39)$$

Thus, by assumption 22.33 and triangle inequality,

$$\|S_\lambda(t)\| \leq e^{-\lambda t} \sum_{k=0}^{\infty} \frac{t^k}{k!} \lambda^{2k} \frac{M}{\lambda^k} = e^{-\lambda t} M e^{\lambda t} = M \quad (22.40)$$

Next is an easy observation about operators commuting.

Lemma 22.8.10 For $\lambda, \mu > 0$, $(\lambda I - A)^{-1}$ and $(\mu I - A)^{-1}$ commute.

Proof: Suppose

$$y = (\mu I - A)^{-1} (\lambda I - A)^{-1} x \quad (22.41)$$

$$z = (\lambda I - A)^{-1} (\mu I - A)^{-1} x \quad (22.42)$$

I need to show $y = z$. This follows from the observation that

$$(\lambda I - A)(\mu I - A)y = (\mu I - A)(\lambda I - A)y = (\mu I - A)(\lambda I - A)z = x$$

■

It follows from the description of $S_\lambda(t)$ in terms of a power series that $S_\lambda(t)$ and $S_\mu(s)$ commute and also A_λ commutes with $S_\mu(t)$ for any t . Indeed, the absolute convergence of the series 22.39 means we can use the Cauchy product to compute the product of these two series and see $S_\lambda(t), S_\mu(t)$ commute. One could also exploit uniqueness and the theory of ordinary differential equations to verify this. I will use this fact in what follows whenever needed.

I want to show that for each $x \in D(A)$,

$$\lim_{\lambda \rightarrow \infty} S_\lambda(t)x \equiv S(t)x$$

where $S(t)$ is the desired semigroup. Let $x \in D(A)$. Then

$$\begin{aligned} S_\mu(t)x - S_\lambda(t)x &= \int_0^t \frac{d}{dr} (S_\lambda(t-r)S_\mu(r)) x dr \\ &= \int_0^t (-S'_\lambda(t-r)S_\mu(r) + S_\lambda(t-r)S'_\mu(r)) x dr \\ &= \int_0^t (S_\lambda(t-r)S_\mu(r)A_\lambda - S_\mu(r)S_\lambda(t-r)A_\mu) x dr \\ &= \int_0^t S_\lambda(t-r)S_\mu(r)(A_\mu x - A_\lambda x) dr \end{aligned}$$

It follows that

$$\|S_\mu(t)x - S_\lambda(t)x\| \leq \int_0^t \|S_\lambda(t-r)S_\mu(r)(A_\mu x - A_\lambda x)\| dr$$

$$\leq M^2 t \|A_\mu x - A_\lambda x\| \leq M^2 t (\|A_\mu x - Ax\| + \|Ax - A_\lambda x\|)$$

Now by Lemma 22.8.9, the right side converges uniformly to 0 in $t \in [0, T]$ an arbitrary finite interval. Denote that to which it converges $S(t)x$. Therefore, $t \rightarrow S(t)x$ is continuous for each $x \in D(A)$ and also from 22.40,

$$\|S(t)x\| = \lim_{\lambda \rightarrow \infty} \|S_\lambda(t)x\| \leq M \|x\|$$

so that $S(t)$ can be extended uniquely to a continuous linear map, still called $S(t)$ defined on all of X which also satisfies $\|S(t)\| \leq M$ since $D(A)$ is dense in X . The uniform convergence on $[0, T]$ implies $t \rightarrow S(t)$ is continuous.

It remains to verify that A generates $S(t)$ and for all x , $\lim_{t \rightarrow 0+} S(t)x - x = 0$. From the above,

$$S_\lambda(t)x = x + \int_0^t S_\lambda(s)A_\lambda x ds \quad (22.43)$$

and so $\lim_{t \rightarrow 0+} \|S_\lambda(t)x - x\| = 0$. By the uniform convergence just shown, there exists λ large enough that for all $t \in [0, \delta]$, $\|S(t)x - S_\lambda(t)x\| < \varepsilon$. Then

$$\begin{aligned} \limsup_{t \rightarrow 0+} \|S(t)x - x\| &\leq \limsup_{t \rightarrow 0+} (\|S(t)x - S_\lambda(t)x\| + \|S_\lambda(t)x - x\|) \\ &\leq \limsup_{t \rightarrow 0+} (\varepsilon + \|S_\lambda(t)x - x\|) \leq \varepsilon \end{aligned}$$

It follows $\lim_{t \rightarrow 0+} S(t)x = x$ because ε is arbitrary.

Next, $\lim_{\lambda \rightarrow \infty} A_\lambda x = Ax$ for all $x \in D(A)$ by Lemma 22.8.9. Therefore, passing to the limit in 22.43 yields from the uniform convergence

$$S(t)x = x + \int_0^t S(s)Ax ds$$

and by continuity of $s \rightarrow S(s)Ax$, it follows

$$\lim_{h \rightarrow 0+} \frac{S(h)x - x}{h} = \lim_{h \rightarrow 0+} \frac{1}{h} \int_0^h S(s)Ax ds = Ax$$

Thus letting B denote the generator of $S(t)$, $D(A) \subseteq D(B)$ and $A = B$ on $D(A)$. It only remains to verify $D(A) = D(B)$.

To do this, let $\lambda > 0$ and consider the following where $y \in X$ is arbitrary.

$$(\lambda I - B)^{-1}y = (\lambda I - B)^{-1}((\lambda I - A)(\lambda I - A)^{-1}y)$$

Now $(\lambda I - A)^{-1}y \in D(A) \subseteq D(B)$ and $A = B$ on $D(A)$ and so

$$(\lambda I - A)(\lambda I - A)^{-1}y = (\lambda I - B)(\lambda I - A)^{-1}y$$

which implies,

$$\begin{aligned} (\lambda I - B)^{-1}y &= \\ (\lambda I - B)^{-1}((\lambda I - B)(\lambda I - A)^{-1}y) &= (\lambda I - A)^{-1}y \end{aligned}$$

Recall from Proposition 22.8.5, an arbitrary element of $D(B)$ is of the form $(\lambda I - B)^{-1}y$ and this has shown every such vector is in $D(A)$, in fact it equals $(\lambda I - A)^{-1}y$. Hence $D(B) \subseteq D(A)$ which shows A generates $S(t)$ and this proves the first half of the theorem.

Next suppose A is the generator of a semigroup $S(t)$ having $\|S(t)\| \leq M$. Then by Proposition 22.8.5 for all $\lambda > 0$, $(\lambda I - A)$ is onto and $(\lambda I - A)^{-1} = \int_0^\infty e^{-\lambda t} S(t) dt$. Thus,

$$\begin{aligned} & \left\| \left((\lambda I - A)^{-1} \right)^n \right\| \\ &= \left\| \int_0^\infty \cdots \int_0^\infty e^{-\lambda(t_1 + \cdots + t_n)} S(t_1 + \cdots + t_n) dt_1 \cdots dt_n \right\| \\ &\leq \int_0^\infty \cdots \int_0^\infty e^{-\lambda(t_1 + \cdots + t_n)} M dt_1 \cdots dt_n = \frac{M}{\lambda^n}. \blacksquare \end{aligned}$$

22.8.3 An Evolution Equation

When Λ generates a continuous semigroup, one can consider a very interesting theorem about evolution equations of the form $y' - \Lambda y = g(t)$ provided $t \rightarrow g(t)$ is C^1 .

Theorem 22.8.11 *Let Λ be the generator of $S(t)$, a continuous semigroup on X , a Banach space and let $t \rightarrow g(t)$ be in $C^1(0, \infty; X)$. Then there exists a unique solution to the initial value problem $y' = \Lambda y + g$, $y(0) = y_0 \in D(\Lambda)$ and it is given by*

$$y(t) = S(t)y_0 + \int_0^t S(t-s)g(s)ds. \quad (22.44)$$

This solution is continuous having continuous derivative and has values in $D(\Lambda)$.

Proof: First I show the following claim.

Claim: For $t > 0$, $\int_0^t S(t-s)g(s)ds \in D(\Lambda)$ and

$$\Lambda \left(\int_0^t S(t-s)g(s)ds \right) = S(t)g(0) - g(t) + \int_0^t S(t-s)g'(s)ds$$

Proof of the claim:

$$\begin{aligned} & \frac{1}{h} \left(S(h) \int_0^t S(t-s)g(s)ds - \int_0^t S(t-s)g(s)ds \right) \\ &= \frac{1}{h} \left(\int_0^t S(t-s+h)g(s)ds - \int_0^t S(t-s)g(s)ds \right) \\ &= \frac{1}{h} \left(\int_{-h}^{t-h} S(t-s)g(s+h)ds - \int_0^t S(t-s)g(s)ds \right) \\ &= \frac{1}{h} \int_{-h}^0 S(t-s)g(s+h)ds + \int_0^{t-h} S(t-s) \frac{g(s+h) - g(s)}{h} ds \\ &\quad - \frac{1}{h} \int_{t-h}^t S(t-s)g(s)ds \end{aligned}$$

Using the estimate in Theorem 22.8.3 on Page 603, the triangle inequality and the uniform convergence of the integrands, the limit as $h \rightarrow 0$ of the above equals

$$S(t)g(0) - g(t) + \int_0^t S(t-s)g'(s)ds$$

which proves the claim since the limit exists and is therefore, $\Lambda \left(\int_0^t S(t-s) g(s) ds \right)$.

Since $y_0 \in D(\Lambda)$,

$$\begin{aligned} S(t) \Lambda y_0 &= S(t) \lim_{h \rightarrow 0} \frac{S(h) y_0 - y_0}{h} = \lim_{h \rightarrow 0} \frac{S(t+h) - S(t)}{h} y_0 \\ &= \lim_{h \rightarrow 0+} \frac{S(h) S(t) y_0 - S(t) y_0}{h} \equiv \Lambda S(t) y_0 \end{aligned} \quad (22.45)$$

This is because the limit exists and so it is by definition the right side. So $S(t) y_0 \in D(\Lambda)$.

Now consider 22.44.

$$\begin{aligned} \frac{y(t+h) - y(t)}{h} &= \frac{S(t+h) - S(t)}{h} y_0 + \\ &\quad \frac{1}{h} \left(\int_0^{t+h} S(t-s+h) g(s) ds - \int_0^t S(t-s) g(s) ds \right) \\ &= \frac{S(t+h) - S(t)}{h} y_0 + \frac{1}{h} \int_t^{t+h} S(t-s+h) g(s) ds \\ &\quad + \frac{1}{h} \left(S(h) \int_0^t S(t-s) g(s) ds - \int_0^t S(t-s) g(s) ds \right) \end{aligned}$$

From the claim and 22.45, the limit of the right side is

$$\begin{aligned} &\Lambda S(t) y_0 + g(t) + \Lambda \left(\int_0^t S(t-s) g(s) ds \right) \\ &= \Lambda \left(S(t) y_0 + \int_0^t S(t-s) g(s) ds \right) + g(t) \end{aligned}$$

Hence $y'(t) = \Lambda y(t) + g(t)$ and from the formula, y' is continuous since by the claim and 22.45 it also equals

$$S(t) \Lambda y_0 + g(t) + S(t) g(0) - g(t) + \int_0^t S(t-s) g'(s) ds$$

which is continuous. The claim and 22.45 also shows $y(t) \in D(\Lambda)$. This proves the existence part of the lemma.

It remains to prove the uniqueness part. It suffices to show that if

$$y' - \Lambda y = 0, \quad y(0) = 0$$

and y is C^1 having values in $D(\Lambda)$, then $y = 0$. Suppose then that y is this way. Letting $0 < s < t$,

$$\begin{aligned} &\frac{d}{ds} (S(t-s) y(s)) \equiv \\ &\lim_{h \rightarrow 0} S(t-s-h) \frac{y(s+h) - y(s)}{h} - \frac{S(t-s) y(s) - S(t-s-h) y(s)}{h} \end{aligned}$$

provided the limit exists. Since y' exists and $y(s) \in D(\Lambda)$, this equals

$$S(t-s) y'(s) - S(t-s) \Lambda y(s) = 0.$$

Let $y^* \in X'$. This has shown that on the open interval $(0, t)$, $s \rightarrow y^*(S(t-s)y(s))$ has a derivative equal to 0. Also from continuity of S and y , this function is continuous on $[0, t]$. Therefore, it is constant on $[0, t]$ by the mean value theorem. At $s=0$, this function equals 0. Therefore, it equals 0 on $[0, t]$. Thus for fixed $s > 0$ and letting $t > s$, $y^*(S(t-s)y(s)) = 0$. Now let t decrease toward s . Then $y^*(y(s)) = 0$ and since y^* was arbitrary, it follows $y(s) = 0$. ■

22.8.4 Adjoints for Closed Operators, Hilbert Space

In Hilbert space, there are some special things which are true.

Definition 22.8.12 *Let A be a densely defined closed operator on H a real Hilbert space. Then A^* is defined as follows.*

$$D(A^*) \equiv \{y \in H : |(Ax, y)| \leq C|x|\}$$

*Then since $D(A)$ is dense, there exists a unique element of H denoted by A^*y such that*

$$(Ax, y) = (x, A^*y)$$

for all $x \in D(A)$.

Lemma 22.8.13 *Let A be closed and densely defined on $D(H) \subseteq H$, a Hilbert space. Then A^* is also closed and densely defined. Also $(A^*)^* = A$. If $(\lambda I - A)^{-1} \in \mathcal{L}(H, H)$, then $(\lambda I - A^*)^{-1} \in \mathcal{L}(H, H)$ and $\left((\lambda I - A)^{-1}\right)^n)^* = (\lambda I - A^*)^{-1})^n$.*

Proof: Denote by $[x, y]$ an ordered pair in $H \times H$. Define $\tau : H \times H \rightarrow H \times H$ by

$$\tau[x, y] \equiv [-y, x]$$

Then the definition of adjoint implies that for $\mathcal{G}(B)$ equal to the graph of B ,

$$\mathcal{G}(A^*) = (\tau\mathcal{G}(A))^\perp \quad (22.46)$$

In this notation the inner product on $H \times H$ with respect to which \perp is defined is given by

$$([x, y], [a, b]) \equiv (x, a) + (y, b).$$

Here is why this is so. For $[x, A^*x] \in \mathcal{G}(A^*)$ it follows that for all $y \in D(A)$

$$([x, A^*x], [-Ay, y]) = -(Ay, x) + (y, A^*x) = 0$$

and so $[x, A^*x] \in (\tau\mathcal{G}(A))^\perp$ which shows $\mathcal{G}(A^*) \subseteq (\tau\mathcal{G}(A))^\perp$. To obtain the other inclusion, let $[a, b] \in (\tau\mathcal{G}(A))^\perp$. This means that for all $x \in D(A)$,

$$([a, b], [-Ax, x]) = 0.$$

In other words, for all $x \in D(A)$, $(Ax, a) = (x, b)$ and so $|(Ax, a)| \leq C|x|$ for all $x \in D(A)$ which shows $a \in D(A^*)$ and $(x, A^*a) = (x, b)$ for all $x \in D(A)$. Therefore, since $D(A)$ is dense, it follows $b = A^*a$ and so $[a, b] \in \mathcal{G}(A^*)$. This shows the other inclusion.

Note that if V is any subspace of the Hilbert space $H \times H$, $(V^\perp)^\perp = \overline{V}$ and S^\perp is always a closed subspace. Also τ and \perp commute. The reason for this is that $[x, y] \in (\tau V)^\perp$ means that $(x, -b) + (y, a) = 0$ for all $[a, b] \in V$ and $[x, y] \in \tau(V^\perp)$ means $[-y, x] \in V^\perp$ so for all $[a, b] \in V$, $(-y, a) + (x, b) = 0$ which says the same thing. It is also clear that $\tau \circ \tau$ has the effect of multiplication by -1 .

It follows from the above description of the graph of A^* that even if $\mathcal{G}(A)$ were not closed it would still be the case that $\mathcal{G}(A^*)$ is closed.

Why is $D(A^*)$ dense? Suppose $z \in D(A^*)^\perp$. Then for all $y \in D(A^*)$ so that $[y, Ay] \in \mathcal{G}(A^*)$, it follows $[z, 0] \in \mathcal{G}(A^*)^\perp = ((\tau\mathcal{G}(A))^\perp)^\perp = \tau\mathcal{G}(A)$ but this implies $[0, z] \in -\mathcal{G}(A)$ and so $z = -A0 = 0$. Thus $D(A^*)$ must be dense since there is no nonzero vector in $D(A^*)^\perp$.

Since A is a closed operator, meaning $\mathcal{G}(A)$ is closed in $H \times H$, it follows from the above formula that

$$\begin{aligned} \mathcal{G}((A^*)^*) &= \left(\tau \left((\tau\mathcal{G}(A))^\perp \right) \right)^\perp = \left(\tau(\tau\mathcal{G}(A))^\perp \right)^\perp \\ &= \left((-\mathcal{G}(A))^\perp \right)^\perp = \left(\mathcal{G}(A)^\perp \right)^\perp = \mathcal{G}(A) \end{aligned}$$

and so $(A^*)^* = A$.

Now consider the final claim. First let $y \in D(A^*) = D(\lambda I - A^*)$. Then letting $x \in H$ be arbitrary,

$$\begin{aligned} &\left(x, \left((\lambda I - A)(\lambda I - A)^{-1} \right)^* y \right) \\ &= \left((\lambda I - A)(\lambda I - A)^{-1} x, y \right) = \left(x, \left((\lambda I - A)^{-1} \right)^* (\lambda I - A^*) y \right) \end{aligned}$$

Thus

$$\left((\lambda I - A)(\lambda I - A)^{-1} \right)^* = I = \left((\lambda I - A)^{-1} \right)^* (\lambda I - A^*) \quad (22.47)$$

on $D(A^*)$. Next let $x \in D(A) = D(\lambda I - A)$ and $y \in H$ arbitrary.

$$(x, y) = \left((\lambda I - A)^{-1} (\lambda I - A)x, y \right) = \left((\lambda I - A)x, \left((\lambda I - A)^{-1} \right)^* y \right)$$

Now it follows $\left| \left((\lambda I - A)x, \left((\lambda I - A)^{-1} \right)^* y \right) \right| \leq |y| |x|$ for any $x \in D(A)$ and so

$$\left((\lambda I - A)^{-1} \right)^* y \in D(A^*)$$

Hence

$$(x, y) = \left(x, (\lambda I - A^*) \left((\lambda I - A)^{-1} \right)^* y \right).$$

Since $x \in D(A)$ is arbitrary and $D(A)$ is dense, it follows

$$(\lambda I - A^*) \left((\lambda I - A)^{-1} \right)^* = I \quad (22.48)$$

From 22.47 and 22.48 it follows $(\lambda I - A^*)^{-1} = \left((\lambda I - A)^{-1} \right)^*$ and $(\lambda I - A^*)$ is one to one and onto with continuous inverse. Finally, from the above,

$$\left((\lambda I - A^*)^{-1} \right)^n = \left(\left((\lambda I - A)^{-1} \right)^* \right)^n = \left(\left((\lambda I - A)^{-1} \right)^n \right)^*. \blacksquare$$

With this preparation, here is an interesting result about the adjoint of the generator of a continuous bounded semigroup. I found this in Balakrishnan [5].

Theorem 22.8.14 *Suppose A is a densely defined closed operator which generates a continuous semigroup, $S(t)$. Then A^* is also a closed densely defined operator which generates $S^*(t)$ and $S^*(t)$ is also a continuous semigroup.*

Proof: First suppose $S(t)$ is also a bounded semigroup, $\|S(t)\| \leq M$. From Lemma 22.8.13 A^* is closed and densely defined. It follows from the Hille Yosida theorem, Theorem 22.8.8 that

$$\left| \left((\lambda I - A)^{-1} \right)^n \right| \leq \frac{M}{\lambda^n}$$

From Lemma 22.8.13 and the fact the adjoint of a bounded linear operator preserves the norm,

$$\begin{aligned} \frac{M}{\lambda^n} &\geq \left| \left(\left((\lambda I - A)^{-1} \right)^n \right)^* \right| = \left| \left((\lambda I - A)^{-1} \right)^* \right|^n \\ &= \left| \left((\lambda I - A^*)^{-1} \right)^n \right| \end{aligned}$$

and so by Theorem 22.8.8 again it follows A^* generates a continuous semigroup, $T(t)$ which satisfies $\|T(t)\| \leq M$. I need to identify $T(t)$ with $S^*(t)$. However, from the proof of Theorem 22.8.8 and Lemma 22.8.13, it follows that for $x \in D(A^*)$ and a suitable sequence $\{\lambda_n\}$,

$$\begin{aligned} (T(t)x, y) &= \left(\lim_{n \rightarrow \infty} e^{-\lambda_n t} \sum_{k=0}^{\infty} \frac{t^k \left(\lambda_n^2 (\lambda_n I - A^*)^{-1} \right)^k}{k!} x, y \right) \\ &= \lim_{n \rightarrow \infty} \left(e^{-\lambda_n t} \sum_{k=0}^{\infty} \frac{t^k \left(\left(\lambda_n^2 (\lambda_n I - A)^{-1} \right)^k \right)^*}{k!} x, y \right) \\ &= \lim_{n \rightarrow \infty} \left(x, e^{-\lambda_n t} \left(\sum_{k=0}^{\infty} \frac{t^k \left(\lambda_n^2 (\lambda_n I - A)^{-1} \right)^k}{k!} \right) y \right) \\ &= (x, S(t)y) = (S^*(t)x, y). \end{aligned}$$

Therefore, since y is arbitrary, $S^*(t) = T(t)$ on $x \in D(A^*)$ a dense set and this shows the two are equal. This proves the proposition in the case where $S(t)$ is also bounded.

Next only assume $S(t)$ is a continuous semigroup. Then by Proposition 22.8.5 there exists $\alpha > 0$ such that

$$\|S(t)\| \leq M e^{\alpha t}.$$

Then consider the operator $-\alpha I + A$ and the bounded semigroup $e^{-\alpha t} S(t)$. For $x \in D(A)$

$$\begin{aligned} \lim_{h \rightarrow 0+} \frac{e^{-\alpha h} S(h)x - x}{h} &= \lim_{h \rightarrow 0+} \left(e^{-\alpha h} \frac{S(h)x - x}{h} + \frac{e^{-\alpha h} - 1}{h} x \right) \\ &= -\alpha x + Ax \end{aligned}$$

Thus $-\alpha I + A$ generates $e^{-\alpha t} S(t)$ and it follows from the first part that $-\alpha I + A^*$ generates $e^{-\alpha t} S^*(t)$. Thus

$$\begin{aligned} -\alpha x + A^* x &= \lim_{h \rightarrow 0+} \frac{e^{-\alpha h} S^*(h) x - x}{h} \\ &= \lim_{h \rightarrow 0+} \left(e^{-\alpha h} \frac{S^*(h) x - x}{h} + \frac{e^{-\alpha h} - 1}{h} x \right) \\ &= -\alpha x + \lim_{h \rightarrow 0+} \frac{S^*(h) x - x}{h} \end{aligned}$$

showing that A^* generates $S^*(t)$. It follows from Proposition 22.8.5 that A^* is closed and densely defined. It is obvious $S^*(t)$ is a semigroup. Why is it continuous? This also follows from the first part of the argument which establishes that $e^{-\alpha t} S^*(t)$ is continuous. This proves the theorem.

22.8.5 Adjoints, Reflexive Banach Space

Here the adjoint of a generator of a semigroup is considered. I will show that the adjoint of the generator generates the adjoint of the semigroup in a reflexive Banach space. This is about as far as you can go although a general but less satisfactory result is given in Yosida [60].

Definition 22.8.15 Let A be a densely defined closed operator on H a real Banach space. Then A^* is defined as follows.

$$D(A^*) \equiv \{y^* \in H' : |y^*(Ax)| \leq C \|x\| \text{ for all } x \in D(A)\}$$

Then since $D(A)$ is dense, there exists a unique element of H' denoted by A^*y such that

$$A^*(y^*)(x) = y^*(Ax)$$

for all $x \in D(A)$.

Lemma 22.8.16 Let A be closed and densely defined on $D(A) \subseteq H$, a Banach space. Then A^* is also closed and densely defined. Also $(A^*)^* = A$. In addition to this, if

$$(\lambda I - A)^{-1} \in \mathcal{L}(H, H),$$

then $(\lambda I - A^*)^{-1} \in \mathcal{L}(H', H')$ and

$$\left(\left((\lambda I - A)^{-1} \right)^n \right)^* = \left((\lambda I - A^*)^{-1} \right)^n$$

Proof: Denote by $[x, y]$ an ordered pair in $H \times H$. Define $\tau : H \times H \rightarrow H \times H$ by

$$\tau[x, y] \equiv [-y, x]$$

A similar notation will apply to $H' \times H'$. Then the definition of adjoint implies that for $\mathcal{G}(B)$ equal to the graph of B ,

$$\mathcal{G}(A^*) = (\tau \mathcal{G}(A))^\perp \quad (22.49)$$

For $S \subseteq H \times H$, define S^\perp by

$$\{[a^*, b^*] \in H' \times H' : a^*(x) + b^*(y) = 0 \text{ for all } [x, y] \in S\}$$

If $S \subseteq H' \times H'$ a similar definition holds.

$$\{[x, y] \in H \times H : a^*(x) + b^*(y) = 0 \text{ for all } [a^*, b^*] \in S\}$$

Here is why 22.49 is so. For $[x^*, A^*x^*] \in \mathcal{G}(A^*)$ it follows that for all $y \in D(A)$

$$x^*(Ay) = A^*x^*(y)$$

and so for all $[y, Ay] \in \mathcal{G}(A)$,

$$-x^*(Ay) + A^*x^*(y) = 0$$

which is what it means to say $[x^*, A^*x^*] \in (\tau\mathcal{G}(A))^\perp$. This shows

$$\mathcal{G}(A^*) \subseteq (\tau\mathcal{G}(A))^\perp$$

To obtain the other inclusion, let $[a^*, b^*] \in (\tau\mathcal{G}(A))^\perp$. This means that for all $[x, Ax] \in \mathcal{G}(A)$,

$$-a^*(Ax) + b^*(x) = 0$$

In other words, for all $x \in D(A)$,

$$|a^*(Ax)| \leq \|b^*\| \|x\|$$

which means by definition, $a^* \in D(A^*)$ and $A^*a^* = b^*$. Thus $[a^*, b^*] \in \mathcal{G}(A^*)$. This shows the other inclusion.

Note that if V is any subspace of $H \times H$, $(V^\perp)^\perp = \overline{V}$. and S^\perp is always a closed subspace. Also τ and \perp commute. The reason for this is that $[x^*, y^*] \in (\tau V)^\perp$ means that

$$-x^*(b) + y^*(a) = 0$$

for all $[a, b] \in V$ and $[x^*, y^*] \in \tau(V^\perp)$ means $[-y^*, x^*] \in -(V^\perp) = V^\perp$ so for all $[a, b] \in V$,

$$-y^*(a) + x^*(b) = 0$$

which says the same thing. It is also clear that $\tau \circ \tau$ has the effect of multiplication by -1 . If $V \subseteq H' \times H'$, the argument for commuting \perp and τ is similar.

It follows from the above description of the graph of A^* that even if $\mathcal{G}(A)$ were not closed it would still be the case that $\mathcal{G}(A^*)$ is closed.

Why is $D(A^*)$ dense? If it is not dense, then by a typical application of the Hahn Banach theorem, there exists $y^{**} \in H''$ such that $y^{**}(D(A^*)) = 0$ but $y^{**} \neq 0$. Since H is reflexive, there exists $y \in H$ such that $x^*(y) = 0$ for all $x^* \in D(A^*)$. Thus

$$[y, 0] \in \mathcal{G}(A^*)^\perp = \left((\tau\mathcal{G}(A))^\perp \right)^\perp = \tau\mathcal{G}(A)$$

and so $[0, y] \in \mathcal{G}(A)$ which means $y = A0 = 0$, a contradiction. Thus $D(A^*)$ is indeed dense. Note this is where it was important to assume the space is reflexive. If you consider

$C([0, 1])$ it is not dense in $L^\infty([0, 1])$ but if $f \in L^1([0, 1])$ satisfies $\int_0^1 f g dm = 0$ for all $g \in C([0, 1])$, then $f = 0$. Hence there is no nonzero $f \in C([0, 1])^\perp$.

Since A is a closed operator, meaning $\mathcal{G}(A)$ is closed in $H \times H$, it follows from the above formula that

$$\begin{aligned} \mathcal{G}((A^*)^*) &= \left(\tau \left((\tau \mathcal{G}(A))^\perp \right) \right)^\perp = \left(\tau (\tau \mathcal{G}(A))^\perp \right)^\perp \\ &= \left((-\mathcal{G}(A))^\perp \right)^\perp = \left(\mathcal{G}(A)^\perp \right)^\perp = \mathcal{G}(A) \end{aligned}$$

and so $(A^*)^* = A$.

Now consider the final claim. First let $y^* \in D(A^*) = D(\lambda I - A^*)$. Then letting $x \in H$ be arbitrary,

$$y^*(x) = \left((\lambda I - A)(\lambda I - A)^{-1} \right)^* y^*(x) = y^* \left((\lambda I - A)(\lambda I - A)^{-1} x \right)$$

Since $y^* \in D(A^*)$ and $(\lambda I - A)^{-1} x \in D(A)$, this equals $(\lambda I - A)^* y^* \left((\lambda I - A)^{-1} x \right)$. Now by definition, this equals $\left((\lambda I - A)^{-1} \right)^* (\lambda I - A)^* y^*(x)$. It follows that for $y^* \in D(A^*)$,

$$\left((\lambda I - A)^{-1} \right)^* (\lambda I - A)^* y^* = \left((\lambda I - A)^{-1} \right)^* (\lambda I - A^*) y^* = y^* \quad (22.50)$$

Next let $y^* \in H'$ be arbitrary and $x \in D(A)$

$$\begin{aligned} y^*(x) &= y^* \left((\lambda I - A)^{-1} (\lambda I - A)x \right) = \left((\lambda I - A)^{-1} \right)^* y^* ((\lambda I - A)x) \\ &= (\lambda I - A)^* \left((\lambda I - A)^{-1} \right)^* y^*(x) \end{aligned}$$

In going from the second to the third line, the first line shows $\left((\lambda I - A)^{-1} \right)^* y^* \in D(A^*)$ and so the third line follows. Since $D(A)$ is dense, it follows

$$(\lambda I - A^*) \left((\lambda I - A)^{-1} \right)^* = I \quad (22.51)$$

Then 22.50 and 22.51 show $\lambda I - A^*$ is one to one and onto from $D(A^*)$ to H' and

$$(\lambda I - A^*)^{-1} = \left((\lambda I - A)^{-1} \right)^*.$$

Finally, from the above, $\left((\lambda I - A^*)^{-1} \right)^n = \left(\left((\lambda I - A)^{-1} \right)^* \right)^n = \left(\left((\lambda I - A)^{-1} \right)^n \right)^*$. This proves the lemma.

With this preparation, here is an interesting result about the adjoint of the generator of a continuous bounded semigroup.

Theorem 22.8.17 *Suppose A is a densely defined closed operator which generates a continuous semigroup, $S(t)$. Then A^* is also a closed densely defined operator which generates $S^*(t)$ and $S^*(t)$ is also a continuous semigroup.*

Proof: First suppose $S(t)$ is also a bounded semigroup, $\|S(t)\| \leq M$. From Lemma 22.8.16 A^* is closed and densely defined. It follows from the Hille Yosida theorem, Theorem 22.8.8 that

$$\left\| \left((\lambda I - A)^{-1} \right)^n \right\| \leq \frac{M}{\lambda^n}$$

From Lemma 22.8.16 and the fact the adjoint of a bounded linear operator preserves the norm,

$$\frac{M}{\lambda^n} \geq \left\| \left(\left((\lambda I - A)^{-1} \right)^n \right)^* \right\| = \left\| \left((\lambda I - A)^{-1} \right)^* \right\|^n = \left\| (\lambda I - A^*)^{-1} \right\|^n$$

and so by Theorem 22.8.8 again it follows A^* generates a continuous semigroup, $T(t)$ which satisfies $\|T(t)\| \leq M$. I need to identify $T(t)$ with $S^*(t)$. However, from the proof of Theorem 22.8.8 and Lemma 22.8.16, it follows that for $x^* \in D(A^*)$ and a suitable sequence $\{\lambda_n\}$,

$$\begin{aligned} T(t)x^*(y) &= \lim_{n \rightarrow \infty} e^{-\lambda_n t} \sum_{k=0}^{\infty} \frac{t^k \left(\lambda_n^2 (\lambda_n I - A^*)^{-1} \right)^k}{k!} x^*(y) \\ &= \lim_{n \rightarrow \infty} e^{-\lambda_n t} \sum_{k=0}^{\infty} \frac{t^k \left(\left(\lambda_n^2 (\lambda_n I - A)^{-1} \right)^k \right)^*}{k!} x^*(y) \\ &= \lim_{n \rightarrow \infty} x^* \left(e^{-\lambda_n t} \left(\sum_{k=0}^{\infty} \frac{t^k \left(\left(\lambda_n^2 (\lambda_n I - A)^{-1} \right)^k \right)}{k!} y \right) \right) = x^*(S(t)y) = S^*(t)x^*(y). \end{aligned}$$

Therefore, since y is arbitrary, $S^*(t) = T(t)$ on $x \in D(A^*)$ a dense set and this shows the two are equal. In particular, $S^*(t)$ is a semigroup because $T(t)$ is. This proves the proposition in the case where $S(t)$ is also bounded.

Next only assume $S(t)$ is a continuous semigroup. Then by Proposition 22.8.5 there exists $\alpha > 0$ such that $\|S(t)\| \leq M e^{\alpha t}$. Then consider the operator $-\alpha I + A$ and the bounded semigroup $e^{-\alpha t} S(t)$. For $x \in D(A)$

$$\lim_{h \rightarrow 0+} \frac{e^{-\alpha h} S(h)x - x}{h} = \lim_{h \rightarrow 0+} \left(e^{-\alpha h} \frac{S(h)x - x}{h} + \frac{e^{-\alpha h} - 1}{h} x \right) = -\alpha x + Ax$$

Thus $-\alpha I + A$ generates $e^{-\alpha t} S(t)$ and it follows from the first part that $-\alpha I + A^*$ generates the semigroup $e^{-\alpha t} S^*(t)$. Thus

$$\begin{aligned} -\alpha x + A^* x &= \lim_{h \rightarrow 0+} \frac{e^{-\alpha h} S^*(h)x - x}{h} \\ &= \lim_{h \rightarrow 0+} \left(e^{-\alpha h} \frac{S^*(h)x - x}{h} + \frac{e^{-\alpha h} - 1}{h} x \right) = -\alpha x + \lim_{h \rightarrow 0+} \frac{S^*(h)x - x}{h} \end{aligned}$$

showing that A^* generates $S^*(t)$. It follows from Proposition 22.8.5 that A^* is closed and densely defined. It is obvious $S^*(t)$ is a semigroup. Why is it continuous? This also follows from the first part of the argument which establishes that $t \rightarrow e^{-\alpha t} S^*(t)x$ is continuous. ■

22.9 Exercises

1. For $f, g \in C([0, 1])$ let $(f, g) = \int_0^1 f(x) \overline{g(x)} dx$. Is this an inner product space? Is it a Hilbert space? What does the Cauchy Schwarz inequality say in this context?
2. Let S denote the unit sphere in a Banach space X , $S \equiv \{x \in X : \|x\| = 1\}$. Show that if Y is a Banach space, then $A \in \mathcal{L}(X, Y)$ is compact if and only if $A(S)$ is precompact, $\overline{A(S)}$ is compact. $A \in \mathcal{L}(X, Y)$ is said to be compact if whenever B is a bounded subset of X , it follows $A(B)$ is a compact subset of Y . In words, A takes bounded sets to precompact sets.
3. \uparrow Show that $A \in \mathcal{L}(X, Y)$ is compact if and only if A^* is compact. **Hint:** Use the result of 2 and the Ascoli Arzela theorem to argue that for S^* the unit ball in X' , there is a subsequence, $\{y_n^*\} \subseteq S^*$ such that y_n^* converges uniformly on the compact set, $\overline{A(S)}$. Thus $\{A^* y_n^*\}$ is a Cauchy sequence in X' . To get the other implication, apply the result just obtained for the operators A^* and A^{**} . Then use results about the embedding of a Banach space into its double dual space.
4. Prove the parallelogram identity, $|x+y|^2 + |x-y|^2 = 2|x|^2 + 2|y|^2$. Next suppose $(X, \|\cdot\|)$ is a real normed linear space and the parallelogram identity holds. Can it be concluded there exists an inner product (\cdot, \cdot) such that $\|x\| = (x, x)^{1/2}$?
5. Let K be a closed, bounded and convex set in \mathbb{R}^n and let $f : K \rightarrow \mathbb{R}^n$ be continuous and let $y \in \mathbb{R}^n$. Show using the Brouwer fixed point theorem there exists a point $x \in K$ such that $P(y - f(x) + x) = x$. Next show that $(y - f(x), z - x) \leq 0$ for all $z \in K$. The existence of this x is known as Browder's lemma and it has great significance in the study of certain types of nonlinear operators. Now suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is continuous and satisfies $\lim_{|x| \rightarrow \infty} \frac{(f(x), x)}{|x|} = \infty$. Show using Browder's lemma that f is onto.
6. Show that every inner product space is uniformly convex. This means that if x_n, y_n are vectors whose norms are no larger than 1 and if $\|x_n + y_n\| \rightarrow 2$, then $\|x_n - y_n\| \rightarrow 0$. More precisely, for every $\varepsilon > 0$, there is a $\delta > 0$ such that if $\|x + y\| > 2 - \delta$ for $\|x\|, \|y\|$ both 1, then $\|x - y\| < \varepsilon$.
7. Let H be separable and let S be an orthonormal set. Show S is countable. **Hint:** How far apart are two elements of the orthonormal set?
8. Suppose $\{x_1, \dots, x_m\}$ is a linearly independent set of vectors in a normed linear space. Show $\text{span}(x_1, \dots, x_m)$ is a closed subspace. Also show every orthonormal set of vectors is linearly independent.
9. Show every Hilbert space, separable or not, has a maximal orthonormal set of vectors.
10. \uparrow Prove Bessel's inequality, which says that if $\{x_n\}_{n=1}^\infty$ is an orthonormal set in H , then for all $x \in H$, $\|x\|^2 \geq \sum_{k=1}^\infty |(x, x_k)|^2$. **Hint:** Show that if $M = \text{span}(x_1, \dots, x_n)$, then $Px = \sum_{k=1}^n x_k(x, x_k)$. Then observe $\|x\|^2 = \|x - Px\|^2 + \|Px\|^2$.
11. \uparrow Show S is a maximal orthonormal set if and only if $\text{span}(S)$ is dense in H , where $\text{span}(S)$ is defined as $\text{span}(S) \equiv \{\text{all finite linear combinations of elements of } S\}$.

12. \uparrow Suppose $\{x_n\}_{n=1}^\infty$ is a maximal orthonormal set. Show that $x = \sum_{n=1}^\infty (x, x_n)x_n \equiv \lim_{N \rightarrow \infty} \sum_{n=1}^N (x, x_n)x_n$ and $\|x\|^2 = \sum_{i=1}^\infty |(x, x_i)|^2$. Show $(x, y) = \sum_{n=1}^\infty (x, x_n)(y, x_n)$. **Hint:** For the last part of this, you might proceed as follows. Show $((x, y)) \equiv \sum_{n=1}^\infty (x, x_n)(y, x_n)$ is a well defined inner product on the Hilbert space which delivers the same norm as the original inner product. Then you could verify that there exists a formula for the inner product in terms of the norm and conclude the two inner products, (\cdot, \cdot) and $((\cdot, \cdot))$ must coincide.
13. Suppose X is an infinite dimensional Banach space and suppose $\{x_1 \cdots x_n\}$ are linearly independent with $\|x_i\| = 1$. By Problem 8 $\text{span}(x_1 \cdots x_n) \equiv X_n$ is a closed linear subspace of X . Now let $z \notin X_n$ and pick $y \in X_n$ such that $\|z - y\| \leq 2 \text{dist}(z, X_n)$ and let $x_{n+1} = \frac{z-y}{\|z-y\|}$. Show the sequence $\{x_k\}$ satisfies $\|x_n - x_k\| \geq 1/2$ whenever $k < n$. Now show the unit ball $\{x \in X : \|x\| \leq 1\}$ in a normed linear space is compact if and only if X is finite dimensional. **Hint:** $\left\| \frac{z-y}{\|z-y\|} - x_k \right\| = \left\| \frac{z-y-x_k\|z-y\|}{\|z-y\|} \right\|$.
14. Show that if A is a self adjoint operator on a Hilbert space and $Ay = \lambda y$ for λ a complex number and $y \neq 0$, then λ must be real. Also verify that if A is self adjoint and $Ax = \mu x$ while $Ay = \lambda y$, then if $\mu \neq \lambda$, it must be the case that $(x, y) = 0$.
15. Theorem 22.8.11 gives the the existence and uniqueness for an evolution equation of the form $y' - \Lambda y = g$, $y(0) = y_0 \in D(\Lambda)$ where g is in $C^1(0, \infty; H)$ for H a Banach space. Recall Λ was the generator of a continuous semigroup $S(h)$. Generalize this to an equation of the form

$$y' - \Lambda y = g + Ly, \quad y(0) = y_0 \in H$$

where L is a continuous linear map. **Hint:** You might consider $\Lambda + L$ and show it generates a continuous semigroup. Then apply the theorem.

16. Generalize Theorem 22.8.11 in case you know that for each $t > 0, S(t)x \in D(\Lambda)$. You might see about removing the differentiability of g as a requirement and maybe the assumption that $y_0 \in D(\Lambda)$. Analytic semigroups have this property. There we typically start with the closed operator and construct the semigroup $S(t)$ using methods from complex analysis.

Chapter 23

Representation Theorems

23.1 Radon Nikodym Theorem

This chapter is on various representation theorems. The first theorem, the Radon Nikodym Theorem, is a representation theorem for one measure in terms of another. This important theorem represents one measure in terms of another. It is Theorem 10.13.7 on Page 302. For a very different approach, see [50] which has a nice proof due to Von Neumann which is based not on the Hahn decomposition of a signed measure, but on the Riesz representation theorem in Hilbert space.

Definition 23.1.1 Let μ and λ be two measures defined on a σ -algebra \mathcal{S} , of subsets of a set, Ω . λ is absolutely continuous with respect to μ , written as $\lambda \ll \mu$, if $\lambda(E) = 0$ whenever $\mu(E) = 0$. A complex measure λ defined on a σ -algebra \mathcal{S} is one which has the property that if the E_i are distinct and measurable, then $\lambda(\cup_i E_i) = \sum_i \lambda(E_i)$. It is a complex measure because each $\lambda(E_i) \in \mathbb{C}$.

Recall Corollary 10.13.11 on Page 302 which involves the case where λ is a signed measure. I am stating it next for convenience.

Corollary 23.1.2 Let μ be a finite measure and λ a signed measure ($\lambda(E) \in \mathbb{R}$) with $\lambda \ll \mu$ meaning that if $\mu(E) = 0$ then $\lambda(E) = 0$. Then there exists $h \in L^1$ such that $\lambda(E) = \int_E h d\mu$.

There is an easy corollary to this which includes complex measures.

Theorem 23.1.3 Let λ be a complex measure and $\lambda \ll \mu$ for μ a finite measure. Then there exists $h \in L^1$ such that $\lambda(E) = \int_E h d\mu$.

Proof: Let $(\operatorname{Re} \lambda)(E) = \operatorname{Re}(\lambda(E))$ with $\operatorname{Im} \lambda$ defined similarly. Then these are signed measures and so there are functions f_1, f_2 in L^1 such that $\operatorname{Re} \lambda(E) = \int_E f_1 d\mu$, $\operatorname{Im} \lambda(E) = \int_E f_2 d\mu$. Then $h \equiv f_1 + if_2$ satisfies the necessary condition. ■

More general versions of the Radon Nikodym theorem available. To see one of these, one can read the treatment in Hewitt and Stromberg [26]. This involves the notion of decomposable measure spaces, a generalization of σ finite.

23.2 Vector Measures

A vector measure is a generalization of a signed or complex measure and it can have values in any topological vector space. Whole books have been written on this subject. See for example the book by Diestel and Uhl [13] titled Vector Measures. I will emphasize only normed linear spaces. This section is about representing one of these measures with respect to its total variation.

Definition 23.2.1 Let $(V, \|\cdot\|)$ be a normed linear space and let (Ω, \mathcal{S}) be a measure space. A function $\mu : \mathcal{S} \rightarrow V$ is a vector measure if μ is countably additive. That is, if $\{E_i\}_{i=1}^\infty$ is a sequence of disjoint sets of \mathcal{S} ,

$$\mu(\cup_{i=1}^\infty E_i) = \sum_{i=1}^\infty \mu(E_i).$$

A signed measure is an example as in Definition 10.13.2 on Page 300.

Note that it makes sense to take finite sums because it is given that μ has values in a vector space in which vectors can be summed. In the above, $\mu(E_i)$ is a vector. It might be a point in \mathbb{R}^p or in any other vector space. In many of the most important applications, it is a vector in some sort of function space which may be infinite dimensional. The infinite sum has the usual meaning. That is $\sum_{i=1}^{\infty} \mu(E_i) = \lim_{n \rightarrow \infty} \sum_{i=1}^n \mu(E_i)$ where the limit takes place relative to the norm on V .

Definition 23.2.2 Let (Ω, \mathcal{S}) be a measure space and let μ be a vector measure defined on \mathcal{S} . A subset, $\pi(E)$, of \mathcal{S} is called a partition of E if $\pi(E)$ consists of finitely many disjoint sets of \mathcal{S} and $\cup \pi(E) = E$. Let

$$|\mu|(E) = \sup \left\{ \sum_{F \in \pi(E)} \|\mu(F)\| : \pi(E) \text{ is a partition of } E \right\}.$$

$|\mu|$ is called the total variation of μ .

The next theorem may seem a little surprising. It states that, if finite, the total variation is a nonnegative measure.

Theorem 23.2.3 If $|\mu|(\Omega) < \infty$, then $|\mu|$ is a measure on \mathcal{S} . Even if $|\mu|(\Omega) = \infty$, $|\mu|(\cup_{i=1}^{\infty} E_i) \leq \sum_{i=1}^{\infty} |\mu|(E_i)$. That is $|\mu|$ is always subadditive and $|\mu|(A) \leq |\mu|(B)$ whenever $A, B \in \mathcal{S}$ with $A \subseteq B$. In earlier terminology, $|\mu|$ is an outer measure.

Proof: Consider the last claim. Let $a < |\mu|(A)$ and let $\pi(A)$ be a partition of A such that $a < \sum_{F \in \pi(A)} \|\mu(F)\|$. Then $\pi(A) \cup \{B \setminus A\}$ is a partition of B and

$$|\mu|(B) \geq \sum_{F \in \pi(A)} \|\mu(F)\| + \|\mu(B \setminus A)\| > a.$$

Since this is true for all such a , it follows $|\mu|(B) \geq |\mu|(A)$ as claimed.

Let $\{E_j\}_{j=1}^{\infty}$ be a sequence of disjoint sets of \mathcal{S} and let $E_{\infty} = \cup_{j=1}^{\infty} E_j$. Then letting $a < |\mu|(E_{\infty})$, it follows from the definition of total variation there exists a partition of E_{∞} , $\pi(E_{\infty}) = \{A_1, \dots, A_n\}$ such that $a < \sum_{i=1}^n \|\mu(A_i)\|$. Also, $A_i = \cup_{j=1}^{\infty} A_i \cap E_j$ and so by the triangle inequality, $\|\mu(A_i)\| \leq \sum_{j=1}^{\infty} \|\mu(A_i \cap E_j)\|$. Therefore, by the above, and either Fubini's theorem or Lemma 2.5.4 on Page 65,

$$a < \sum_{i=1}^n \overbrace{\sum_{j=1}^{\infty} \|\mu(A_i \cap E_j)\|}^{\geq \|\mu(A_i)\|} = \sum_{j=1}^{\infty} \sum_{i=1}^n \|\mu(A_i \cap E_j)\| \leq \sum_{j=1}^{\infty} |\mu|(E_j)$$

because $\{A_i \cap E_j\}_{i=1}^n$ is a partition of E_j .

Since a is arbitrary, this shows $|\mu|(\cup_{j=1}^{\infty} E_j) \leq \sum_{j=1}^{\infty} |\mu|(E_j)$. If the sets, E_j are not disjoint, let $F_1 = E_1$ and if F_n has been chosen, let $F_{n+1} \equiv E_{n+1} \setminus \cup_{i=1}^n E_i$. Thus the sets, F_i are disjoint and $\cup_{i=1}^{\infty} F_i = \cup_{i=1}^{\infty} E_i$. Therefore,

$$|\mu|(\cup_{j=1}^{\infty} E_j) = |\mu|(\cup_{j=1}^{\infty} F_j) \leq \sum_{j=1}^{\infty} |\mu|(F_j) \leq \sum_{j=1}^{\infty} |\mu|(E_j)$$

and proves $|\mu|$ is always subadditive as claimed, regardless of whether $|\mu|(\Omega) < \infty$.

Now suppose $|\mu|(\Omega) < \infty$ and let E_1 and E_2 be sets of \mathcal{S} such that $E_1 \cap E_2 = \emptyset$ and let $\{A_1^i \cdots A_{n_i}^i\} = \pi(E_i)$, a partition of E_i which is chosen such that

$$|\mu|(E_i) - \varepsilon < \sum_{j=1}^{n_i} \|\mu(A_j^i)\| \quad i = 1, 2.$$

Such a partition exists because of the definition of the total variation. Considering the sets which are contained in either of $\pi(E_1)$ or $\pi(E_2)$, it follows this finite collection of sets is a partition of $E_1 \cup E_2$ denoted by $\pi(E_1 \cup E_2)$. Then by the above inequality and the definition of total variation,

$$|\mu|(E_1 \cup E_2) \geq \sum_{F \in \pi(E_1 \cup E_2)} \|\mu(F)\| > |\mu|(E_1) + |\mu|(E_2) - 2\varepsilon,$$

which shows that since $\varepsilon > 0$ was arbitrary, $|\mu|(E_1 \cup E_2) \geq |\mu|(E_1) + |\mu|(E_2)$. By induction, whenever the E_i are disjoint, $|\mu|(\cup_{j=1}^n E_j) \geq \sum_{j=1}^n |\mu|(E_j)$. Therefore,

$$\sum_{j=1}^{\infty} |\mu|(E_j) \geq |\mu|(\cup_{j=1}^{\infty} E_j) \geq |\mu|(\cup_{j=1}^n E_j) \geq \sum_{j=1}^n |\mu|(E_j).$$

Now let $n \rightarrow \infty$. Thus, $|\mu|(\cup_{j=1}^{\infty} E_j) = \sum_{j=1}^{\infty} |\mu|(E_j)$ which shows that $|\mu|$ is a measure as claimed. ■

The following corollary is interesting. It concerns the case that μ is only finitely additive.

Corollary 23.2.4 *Suppose (Ω, \mathcal{F}) is a set with a σ algebra of subsets \mathcal{F} and suppose $\mu : \mathcal{F} \rightarrow \mathbb{C}$ is only finitely additive. That is, $\mu(\cup_{i=1}^n E_i) = \sum_{i=1}^n \mu(E_i)$ whenever the E_i are disjoint. Then $|\mu|$, defined in the same way as above, is also finitely additive provided $|\mu|$ is finite.*

Proof: Say $E \cap F = \emptyset$ for $E, F \in \mathcal{F}$. Let $\pi(E), \pi(F)$ suitable partitions for which the following holds.

$$|\mu|(E \cup F) \geq \sum_{A \in \pi(E)} |\mu(A)| + \sum_{B \in \pi(F)} |\mu(B)| \geq |\mu|(E) + |\mu|(F) - 2\varepsilon.$$

Since ε is arbitrary, $|\mu|(E \cap F) \geq |\mu|(E) + |\mu|(F)$. Similar considerations apply to any finite union of disjoint sets. That is, if the E_i are disjoint, then $|\mu|(\cup_{i=1}^n E_i) \geq \sum_{i=1}^n |\mu|(E_i)$.

Now let $E = \cup_{i=1}^n E_i$ where the E_i are disjoint. Then letting $\pi(E)$ be a suitable partition of E ,

$$|\mu|(E) - \varepsilon \leq \sum_{F \in \pi(E)} |\mu(F)|,$$

it follows that

$$\begin{aligned} |\mu|(E) &\leq \varepsilon + \sum_{F \in \pi(E)} |\mu(F)| = \varepsilon + \sum_{F \in \pi(E)} \left| \sum_{i=1}^n \mu(F \cap E_i) \right| \\ &\leq \varepsilon + \sum_{i=1}^n \sum_{F \in \pi(E)} |\mu(F \cap E_i)| \leq \varepsilon + \sum_{i=1}^n |\mu|(E_i) \end{aligned}$$

Since ε is arbitrary, this shows $|\mu|(\cup_{i=1}^n E_i) \leq \sum_{i=1}^n |\mu|(E_i)$. Thus $|\mu|$ is finitely additive. ■

In the case that λ is a complex measure, it is always the case that $|\lambda|(\Omega) < \infty$. First is a lemma.

Lemma 23.2.5 *Suppose λ is a real valued measure called a signed measure (Definition 10.13.2). Then $|\lambda|$ is a finite measure.*

Proof: Suppose $\lambda : \mathcal{F} \rightarrow \mathbb{R}$ is a vector measure (signed measure by Definition 10.13.2). By the Hahn decomposition, Theorem 10.13.5 on Page 301, $\Omega = P \cup N$ where P is a positive set and N is a negative one. Then on N , $-\lambda$ acts like a measure in the sense that if $A \subseteq B$ and A, B measurable subsets of N , then $-\lambda(A) \leq -\lambda(B)$. Similarly λ is a measure on P .

$$\begin{aligned} \sum_{F \in \pi(\Omega)} |\lambda(F)| &\leq \sum_{F \in \pi(\Omega)} (|\lambda(F \cap P)| + |\lambda(F \cap N)|) \\ &= \sum_{F \in \pi(\Omega)} \lambda(F \cap P) + \sum_{F \in \pi(\Omega)} -\lambda(F \cap N) \\ &= \lambda((\cup_{F \in \pi(\Omega)} F) \cap P) + -\lambda((\cup_{F \in \pi(\Omega)} F) \cap N) \leq \lambda(P) + |\lambda(N)| \end{aligned}$$

It follows that $|\lambda|(\Omega) < \lambda(P) + |\lambda(N)|$ and so $|\lambda|$ has finite total variation. ■

Theorem 23.2.6 *Suppose λ is a complex measure on (Ω, \mathcal{S}) where \mathcal{S} is a σ -algebra of subsets of Ω . Then $|\lambda|(\Omega) < \infty$.*

Proof: If λ is a vector measure with values in \mathbb{C} , $\text{Re } \lambda$ and $\text{Im } \lambda$ have values in \mathbb{R} . Then

$$\begin{aligned} \sum_{F \in \pi(\Omega)} |\lambda(F)| &\leq \sum_{F \in \pi(\Omega)} |\text{Re } \lambda(F)| + |\text{Im } \lambda(F)| \\ &= \sum_{F \in \pi(\Omega)} |\text{Re } \lambda(F)| + \sum_{F \in \pi(\Omega)} |\text{Im } \lambda(F)| \\ &\leq |\text{Re } \lambda|(\Omega) + |\text{Im } \lambda|(\Omega) < \infty \end{aligned}$$

thanks to Lemma 23.2.5. ■

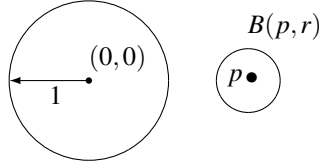
The following corollary is about representing a complex measure λ in terms of its total variation $|\lambda|$. It is like representing a complex number in the form $re^{i\theta}$. In particular, I want to show that in the Radon Nikodym theorem $\left| \frac{d\lambda}{d\mu} \right| = 1$ a.e. First is a lemma which is interesting for its own sake and shows $\left| \frac{d\lambda}{d\mu} \right| \leq 1$.

Lemma 23.2.7 *Suppose $(\Omega, \mathcal{S}, \mu)$ is a measure space and f is a function in $L^1(\Omega, \mu)$ with the property that*

$$\left| \int_E f d\mu \right| \leq \mu(E)$$

for all $E \in \mathcal{S}$. Then $|f| \leq 1$ a.e.

Proof of the lemma: Consider the following picture where $B(p, r) \cap B(0, 1) = \emptyset$.



Let $E = f^{-1}(B(p, r))$. In fact $\mu(E) = 0$. If $\mu(E) \neq 0$ then

$$\left| \frac{1}{\mu(E)} \int_E f d\mu - p \right| = \left| \frac{1}{\mu(E)} \int_E (f - p) d\mu \right| \leq \frac{1}{\mu(E)} \int_E |f - p| d\mu < r$$

because on E , $|f(\omega) - p| < r$. Hence $\frac{1}{\mu(E)} \int_E f d\mu$ is closer to p than r and so

$$\left| \frac{1}{\mu(E)} \int_E f d\mu \right| > 1.$$

Refer to the picture. However, this contradicts the assumption of the lemma. It follows $\mu(E) = 0$. Since the set of complex numbers z such that $|z| > 1$ is an open set, it equals the union of countably many balls, $\{B_i\}_{i=1}^\infty$. Therefore,

$$\mu(f^{-1}(\{z \in \mathbb{C} : |z| > 1\})) = \mu(\cup_{k=1}^\infty f^{-1}(B_k)) \leq \sum_{k=1}^\infty \mu(f^{-1}(B_k)) = 0.$$

Thus $|f(\omega)| \leq 1$ a.e. as claimed. ■

Note that the above argument would work with essentially no change if \mathbb{C} were replaced with V a separable normed linear space.

Corollary 23.2.8 *Let λ be a complex vector measure with $|\lambda|(\Omega) < \infty$.¹ Then there exists a unique $f \in L^1(\Omega)$ such that $\lambda(E) = \int_E f d|\lambda|$. Furthermore, $|f| = 1$ for $|\lambda|$ a.e. This is called the polar decomposition of λ . We write $d\lambda = f d|\lambda|$ sometimes.*

Proof: Letting $\mu = |\lambda|$ in Theorem 23.1.3, the first claim follows because $\lambda \ll |\lambda|$ and so such an L^1 function exists and is unique. It is required to show $|f| = 1$ a.e. If $|\lambda|(E) \neq 0$, $\left| \frac{\lambda(E)}{|\lambda|(E)} \right| = \left| \frac{1}{|\lambda|(E)} \int_E f d|\lambda| \right| \leq 1$. Therefore by Lemma 23.2.7, $|f| \leq 1$, $|\lambda|$ a.e. Now let

$$E_n = \left[|f| \leq 1 - \frac{1}{n} \right].$$

Let $\{F_1, \dots, F_m\}$ be a partition of E_n . Then

$$\begin{aligned} \sum_{i=1}^m |\lambda(F_i)| &= \sum_{i=1}^m \left| \int_{F_i} f d|\lambda| \right| \leq \sum_{i=1}^m \int_{F_i} |f| d|\lambda| \\ &\leq \sum_{i=1}^m \int_{F_i} \left(1 - \frac{1}{n} \right) d|\lambda| = \sum_{i=1}^m \left(1 - \frac{1}{n} \right) |\lambda|(F_i) = |\lambda|(E_n) \left(1 - \frac{1}{n} \right). \end{aligned}$$

Then taking the supremum over all partitions, $|\lambda|(E_n) \leq \left(1 - \frac{1}{n} \right) |\lambda|(E_n)$ which shows $|\lambda|(E_n) = 0$. Hence $|\lambda|([|f| < 1]) = 0$ because $[|f| < 1] = \cup_{n=1}^\infty E_n$. ■

Next is a specific case which leads to complex measures.

¹As proved above, the assumption that $|\lambda|(\Omega) < \infty$ is redundant.

Corollary 23.2.9 Suppose (Ω, \mathcal{S}) is a measure space and μ is a finite nonnegative measure on \mathcal{S} . Then for $h \in L^1(\mu)$, define a complex measure, λ by $\lambda(E) \equiv \int_E h d\mu$. Then $|\lambda|(E) = \int_E |h| d\mu$. Furthermore, $|h| = \bar{g}h$ where $gd|\lambda|$ is the polar decomposition of λ , defined by $\lambda(E) = \int_E gd|\lambda|$.

Proof: From Corollary 23.2.8, there exists g such that $|g| = 1, |\lambda|$ a.e. and for all $E \in \mathcal{S}$

$$\lambda(E) = \int_E gd|\lambda|, \lambda(E) \equiv \int_E h d\mu, \text{ so } \int_E h d\mu = \int_E gd|\lambda| \quad (23.1)$$

Since $|g| = 1$, there is a sequence s_n of simple functions converging pointwise to \bar{g} with $|s_n| \leq 1$. (Approximate the positive and negative parts of the real and imaginary parts of \bar{g} with an increasing sequence of simple functions. Then assemble these to get s_n . See Theorem 10.7.6 on Page 286.) Then from 23.1, $\int_E gs_n d|\lambda| = \int_E s_n h d\mu$. Passing to the limit using the dominated convergence theorem, $\int_E d|\lambda| = \int_E \bar{g}h d\mu$. It follows $\bar{g}h \geq 0$ a.e. and $|\bar{g}| = 1$ a.e. Therefore, $|h| = |\bar{g}h| = \bar{g}h$. It follows from the above, that

$$|\lambda|(E) = \int_E d|\lambda| = \int_E \bar{g}h d\mu = \int_E |h| d\mu \quad \blacksquare$$

Formally: If $d\lambda = h d\mu$, and $d\lambda = gd|\lambda|, |g| = 1$, then $d|\lambda| = |h| d\mu$ and so you might expect to have $d\lambda = g|h| d\mu$ so $g|h| = h$ and so $|h| = \bar{g}h$. That which should be true is. Emphasizing the most significant part of this, if $d\lambda = h d\mu$, then $d|\lambda| = |h| d\mu$.

23.3 Representation for the Dual Space of L^p

Recall the concept of the dual space of a Banach space in the chapter on Banach space starting on Page 533. The next topic deals with the dual space of L^p for $p \geq 1$ in the case where the measure space is σ finite or finite. In what follows $q = \infty$ if $p = 1$ and otherwise, $\frac{1}{p} + \frac{1}{q} = 1$. In what follows, $|\cdot|$ is the usual norm on \mathbb{C} .

Theorem 23.3.1 (Riesz representation theorem) Let $\infty > p > 1$ and let $(\Omega, \mathcal{S}, \mu)$ be a finite measure space. If $\Lambda \in (L^p(\Omega))'$, then there exists a unique $h \in L^q(\Omega)$ such that

$$\Lambda f = \int_{\Omega} h f d\mu.$$

This function satisfies $\|h\|_q = \|\Lambda\|$ where $\|\Lambda\|$ is the operator norm of Λ .

Proof: (Uniqueness) If h_1 and h_2 both represent Λ , consider

$$f = |h_1 - h_2|^{q-2}(\overline{h_1} - \overline{h_2}),$$

where \bar{h} denotes complex conjugation. By Holder's inequality, it is easy to see that $f \in L^p(\Omega)$. Thus $0 = \Lambda f - \Lambda f = \int h_1 |h_1 - h_2|^{q-2}(\overline{h_1} - \overline{h_2}) - h_2 |h_1 - h_2|^{q-2}(\overline{h_1} - \overline{h_2}) d\mu = \int |h_1 - h_2|^q d\mu$. Therefore $h_1 = h_2$ and this proves uniqueness in every case regardless whether μ is finite.

Now let $\lambda(E) = \Lambda(\mathcal{X}_E)$. Since this is a finite measure space, \mathcal{X}_E is an element of $L^p(\Omega)$ and so it makes sense to write $\Lambda(\mathcal{X}_E)$. In fact λ is a complex measure having finite total variation. First I show that it is a measure. Then by Theorem 23.2.6, λ has finite total variation.

If $\{E_i\}_{i=1}^\infty$ is a sequence of disjoint sets of \mathcal{S} , let $F_n = \cup_{i=1}^n E_i$, $F = \cup_{i=1}^\infty E_i$. Then by the Dominated Convergence theorem, $\|\mathcal{X}_{F_n} - \mathcal{X}_F\|_p \rightarrow 0$. Therefore, by continuity of Λ ,

$$\lambda(F) \equiv \Lambda(\mathcal{X}_F) = \lim_{n \rightarrow \infty} \Lambda(\mathcal{X}_{F_n}) = \lim_{n \rightarrow \infty} \sum_{k=1}^n \Lambda(\mathcal{X}_{E_k}) = \sum_{k=1}^\infty \lambda(E_k).$$

This shows λ is a complex measure. Since a similar theorem will be proved in which λ has values in an infinite dimensional space, I will prove this directly without using Theorem 23.2.6. Let A_1, \dots, A_n be a partition of Ω . $|\Lambda(\mathcal{X}_{A_i})| = w_i(\Lambda(\mathcal{X}_{A_i})) = \Lambda(w_i \mathcal{X}_{A_i})$ for some $w_i \in \mathbb{C}$, $|w_i| = 1$. Thus

$$\begin{aligned} \sum_{i=1}^n |\lambda(A_i)| &= \sum_{i=1}^n |\Lambda(\mathcal{X}_{A_i})| = \Lambda\left(\sum_{i=1}^n w_i \mathcal{X}_{A_i}\right) \\ &\leq \|\Lambda\| \left(\int \left|\sum_{i=1}^n w_i \mathcal{X}_{A_i}\right|^p d\mu\right)^{\frac{1}{p}} = \|\Lambda\| \left(\int_{\Omega} d\mu\right)^{\frac{1}{p}} = \|\Lambda\| \mu(\Omega)^{\frac{1}{p}}. \end{aligned}$$

This is because if $x \in \Omega$, x is contained in exactly one of the A_i and so the absolute value of the sum in the first integral above is equal to 1. Therefore $|\lambda|(\Omega) < \infty$ because this was an arbitrary partition. with $|\lambda|$ finite.

It is also clear from the definition of λ that $\lambda \ll \mu$. Therefore, by the Radon Nikodym theorem, there exists $h \in L^1(\Omega)$ with $\lambda(E) = \int_E h d\mu = \Lambda(\mathcal{X}_E)$. Actually $h \in L^q$ and satisfies the other conditions above. This is shown next.

Let $s = \sum_{i=1}^m c_i \mathcal{X}_{E_i}$ be a simple function. Then since Λ is linear,

$$\Lambda(s) = \sum_{i=1}^m c_i \Lambda(\mathcal{X}_{E_i}) = \sum_{i=1}^m c_i \int_{E_i} h d\mu = \int h s d\mu. \quad (23.2)$$

Claim: If f is uniformly bounded and measurable, then

$$\Lambda(f) = \int h f d\mu.$$

Proof of claim: Since f is bounded and measurable, there exists a sequence of simple functions, $\{s_n\}$ which converges to f pointwise and in $L^p(\Omega)$, $|s_n| \leq |f|$. This follows from Theorem 9.1.6 on Page 239 upon breaking f up into positive and negative parts of real and complex parts. In fact this theorem gives uniform convergence. Then

$$\Lambda(f) = \lim_{n \rightarrow \infty} \Lambda(s_n) = \lim_{n \rightarrow \infty} \int h s_n d\mu = \int h f d\mu,$$

the first equality holding because of continuity of Λ , the second following from 23.2 and the third holding by the dominated convergence theorem.

This is a very nice formula but it still has not been shown that $h \in L^q(\Omega)$.

Let $E_n = \{x : |h(x)| \leq n\}$. Thus $|h \mathcal{X}_{E_n}| \leq n$. Then

$$|h \mathcal{X}_{E_n}|^{q-2} (\bar{h} \mathcal{X}_{E_n}) \in L^p(\Omega).$$

By the claim, it follows that

$$\|h \mathcal{X}_{E_n}\|_q^q = \int h |h \mathcal{X}_{E_n}|^{q-2} (\bar{h} \mathcal{X}_{E_n}) d\mu = \Lambda(|h \mathcal{X}_{E_n}|^{q-2} (\bar{h} \mathcal{X}_{E_n}))$$

$$\leq \|\Lambda\| \| |h\mathcal{X}_{E_n}|^{q-2} (\bar{h}\mathcal{X}_{E_n}) \|_p = \left(\int |h\mathcal{X}_{E_n}|^q d\mu \right)^{1/p} = \|\Lambda\| \|h\mathcal{X}_{E_n}\|_q^{\frac{q}{p}},$$

because $q-1 = q/p$ and so it follows that $\|h\mathcal{X}_{E_n}\|_q \leq \|\Lambda\|$. Letting $n \rightarrow \infty$, the monotone convergence theorem implies

$$\|h\|_q \leq \|\Lambda\|. \quad (23.3)$$

Now that h has been shown to be in $L^q(\Omega)$, it follows from 23.2 and the density of the simple functions, Theorem 12.2.1 on Page 362, that $\Lambda f = \int h f d\mu$ for all $f \in L^p(\Omega)$. It only remains to verify the last claim that $\|h\|_q = \|\Lambda\|$ not just 23.3. However, from the definition and Holder's inequality and 23.3, $\|\Lambda\| \equiv \sup\{\int h f : \|f\|_p \leq 1\} \leq \|h\|_q \leq \|\Lambda\|$ ■

To represent elements of the dual space of $L^1(\Omega)$, another Banach space is needed.

Definition 23.3.2 Let $(\Omega, \mathcal{S}, \mu)$ be a measure space. $L^\infty(\Omega)$ is the vector space of measurable functions such that for some $M > 0$, $|f(x)| \leq M$ for all x outside of some set of measure zero ($|f(x)| \leq M$ a.e.). Define $f = g$ when $f(x) = g(x)$ a.e. and $\|f\|_\infty \equiv \inf\{M : |f(x)| \leq M \text{ a.e.}\}$.

Theorem 23.3.3 $L^\infty(\Omega)$ is a Banach space.

Proof: It is clear that $L^\infty(\Omega)$ is a vector space. Is $\|\cdot\|_\infty$ a norm?

Claim: If $f \in L^\infty(\Omega)$, then $|f(x)| \leq \|f\|_\infty$ a.e.

Proof of the claim: $\{x : |f(x)| \geq \|f\|_\infty + n^{-1}\} \equiv E_n$ is a set of measure zero according to the definition of $\|f\|_\infty$. Furthermore, $\{x : |f(x)| > \|f\|_\infty\} = \cup_n E_n$ and so it is also a set of measure zero. This verifies the claim.

Now if $\|f\|_\infty = 0$ it follows that $f(x) = 0$ a.e. Also if

$$f, g \in L^\infty(\Omega)$$

then $|f(x) + g(x)| \leq |f(x)| + |g(x)| \leq \|f\|_\infty + \|g\|_\infty$ a.e. and so $\|f\|_\infty + \|g\|_\infty$ serves as one of the constants, M in the definition of $\|f + g\|_\infty$. Therefore, $\|f + g\|_\infty \leq \|f\|_\infty + \|g\|_\infty$. Next let c be a number. Then $|cf(x)| = |c||f(x)| \leq |c|\|f\|_\infty$ and so $\|cf\|_\infty \leq |c|\|f\|_\infty$. Therefore since c is arbitrary, $\|f\|_\infty = \|c(1/c)f\|_\infty \leq \frac{1}{|c|}\|cf\|_\infty$ which implies $|c|\|f\|_\infty \leq \|cf\|_\infty$. Thus $\|\cdot\|_\infty$ is a norm as claimed.

To verify completeness, let $\{f_n\}$ be a Cauchy sequence in $L^\infty(\Omega)$ and use the above claim to get the existence of a set of measure zero, E_{nm} such that for all $x \notin E_{nm}$,

$$|f_n(x) - f_m(x)| \leq \|f_n - f_m\|_\infty.$$

Let $E = \cup_{n,m} E_{nm}$. Thus $\mu(E) = 0$ and for each $x \notin E$, $\{f_n(x)\}_{n=1}^\infty$ is a Cauchy sequence in \mathbb{C} . Let

$$f(x) = \begin{cases} 0 & \text{if } x \in E \\ \lim_{n \rightarrow \infty} f_n(x) & \text{if } x \notin E \end{cases} = \lim_{n \rightarrow \infty} \mathcal{X}_{E^c}(x) f_n(x).$$

Then f is clearly measurable because it is the limit of measurable functions. If

$$F_n = \{x : |f_n(x)| > \|f_n\|_\infty\}$$

and $F = \cup_{n=1}^\infty F_n$, it follows $\mu(F) = 0$ and that for $x \notin F \cup E$,

$$|f(x)| \leq \liminf_{n \rightarrow \infty} |f_n(x)| \leq \liminf_{n \rightarrow \infty} \|f_n\|_\infty < \infty$$

because $\{\|f_n\|_\infty\}$ is a Cauchy sequence. ($|\|f_n\|_\infty - \|f_m\|_\infty| \leq \|f_n - f_m\|_\infty$ by the triangle inequality.) Thus $f \in L^\infty(\Omega)$. Let n be large enough that whenever $m > n$, $\|f_m - f_n\|_\infty < \varepsilon$. Then, if $x \notin E$,

$$|f(x) - f_n(x)| = \lim_{m \rightarrow \infty} |f_m(x) - f_n(x)| \leq \liminf_{m \rightarrow \infty} \|f_m - f_n\|_\infty < \varepsilon.$$

Hence $\|f - f_n\|_\infty < \varepsilon$ for all n large enough. ■

The next theorem is the Riesz representation theorem for $(L^1(\Omega))'$.

Theorem 23.3.4 (Riesz representation theorem) *Let $(\Omega, \mathcal{S}, \mu)$ be a finite measure space. If $\Lambda \in (L^1(\Omega))'$, then there exists a unique $h \in L^\infty(\Omega)$ such that*

$$\Lambda(f) = \int_{\Omega} hf \, d\mu$$

for all $f \in L^1(\Omega)$. If h is the function in $L^\infty(\Omega)$ representing $\Lambda \in (L^1(\Omega))'$, then $\|h\|_\infty = \|\Lambda\|$.

Proof: Just as in the proof of Theorem 23.3.1, there exists a unique $h \in L^1(\Omega)$ such that for all simple functions s ,

$$\Lambda(s) = \int_{\Omega} hs \, d\mu. \quad (23.4)$$

To show $h \in L^\infty(\Omega)$, let $\varepsilon > 0$ be given and let $E = \{x : |h(x)| \geq \|\Lambda\| + \varepsilon\}$. Let $|k| = 1$ and $hk = |h|$. Since the measure space is finite, $k \in L^1(\Omega)$. As in Theorem 23.3.1 let $\{s_n\}$ be a sequence of simple functions converging to k in $L^1(\Omega)$, and pointwise. It follows from the construction in Theorem 9.1.6 on Page 239 that it can be assumed $|s_n| \leq 1$. Therefore

$$\Lambda(k\mathcal{X}_E) = \lim_{n \rightarrow \infty} \Lambda(s_n\mathcal{X}_E) = \lim_{n \rightarrow \infty} \int_E hs_n \, d\mu = \int_E hkd\mu$$

where the last equality holds by the Dominated Convergence theorem. Therefore,

$$\begin{aligned} \|\Lambda\|\mu(E) &\geq |\Lambda(k\mathcal{X}_E)| = \left| \int_{\Omega} hk\mathcal{X}_E \, d\mu \right| = \int_E |h| \, d\mu \\ &\geq (\|\Lambda\| + \varepsilon)\mu(E). \end{aligned}$$

It follows that $\mu(E) = 0$. Since $\varepsilon > 0$ was arbitrary, $\|\Lambda\| \geq \|h\|_\infty$. Since $h \in L^\infty(\Omega)$, the density of the simple functions in $L^1(\Omega)$ and 23.4 imply

$$\Lambda f = \int_{\Omega} hf \, d\mu, \quad \|\Lambda\| \geq \|h\|_\infty. \quad (23.5)$$

This proves the existence part of the theorem. To verify uniqueness, suppose h_1 and h_2 both represent Λ and let $f \in L^1(\Omega)$ be such that $|f| \leq 1$ and $f(h_1 - h_2) = |h_1 - h_2|$. Then $0 = \Lambda f - \Lambda f = \int (h_1 - h_2)f \, d\mu = \int |h_1 - h_2| \, d\mu$. Thus $h_1 = h_2$. Finally, $\|\Lambda\| = \sup\{|\int hf \, d\mu| : \|f\|_1 \leq 1\} \leq \|h\|_\infty \leq \|\Lambda\|$ by 23.5. ■

Next these results are extended to the σ finite case.

Lemma 23.3.5 *Let $(\Omega, \mathcal{S}, \mu)$ be a measure space and suppose there exists a measurable function, r such that $r(x) > 0$ for all x , there exists M such that $|r(x)| < M$ for all x , and $\int r \, d\mu < \infty$. Then for $\Lambda \in (L^p(\Omega, \mu))'$, $p \geq 1$, there exists $h \in L^q(\Omega, \mu)$, $L^\infty(\Omega, \mu)$ if $p = 1$ such that $\Lambda f = \int hf \, d\mu$. Also $\|h\| = \|\Lambda\|$. ($\|h\| = \|h\|_q$ if $p > 1$, $\|h\|_\infty$ if $p = 1$). Here $\frac{1}{p} + \frac{1}{q} = 1$.*

Proof: Define a new measure $\tilde{\mu}$, according to the rule

$$\tilde{\mu}(E) \equiv \int_E r d\mu. \quad (23.6)$$

Thus $\tilde{\mu}$ is a finite measure on \mathcal{S} . For

$$\Lambda \in (L^p(\mu))', \Lambda(f) = \Lambda\left(r^{1/p} \left(r^{-1/p} f\right)\right) = \tilde{\Lambda}\left(r^{-1/p} f\right)$$

where $\tilde{\Lambda}(g) \equiv \Lambda(r^{1/p} g)$. Now $\tilde{\Lambda}$ is in $L^p(\tilde{\mu})'$ because

$$\begin{aligned} |\tilde{\Lambda}(g)| &\equiv \left| \Lambda\left(r^{1/p} g\right) \right| \leq \|\Lambda\| \left(\int_{\Omega} \left| r^{1/p} g \right|^p d\mu \right)^{1/p} \\ &= \|\Lambda\| \left(\int_{\Omega} |g|^p \overbrace{r d\mu}^{d\tilde{\mu}} \right)^{1/p} = \|\Lambda\| \|g\|_{L^p(\tilde{\mu})} \end{aligned}$$

Therefore, by Theorems 23.3.4 and 23.3.1 there exists a unique $h \in L^q(\tilde{\mu})$ which represents $\tilde{\Lambda}$. Here $q = \infty$ if $p = 1$ and satisfies $1/q + 1/p = 1$ otherwise. Then

$$\Lambda(f) = \tilde{\Lambda}\left(r^{-1/p} f\right) = \int_{\Omega} h f r^{-1/p} r d\mu = \int_{\Omega} f \left(h r^{1/q} \right) d\mu$$

Now $h r^{1/q} \equiv \tilde{h} \in L^q(\mu)$ since $h \in L^q(\tilde{\mu})$. In case $p = 1$, $L^q(\tilde{\mu})$ and $L^q(\mu)$ are exactly the same. In this case you have

$$\Lambda(f) = \tilde{\Lambda}(r^{-1} f) = \int_{\Omega} h f r^{-1} r d\mu = \int_{\Omega} f h d\mu$$

Thus the desired representation holds. Then in any case, $|\Lambda(f)| \leq \|\tilde{h}\|_{L^q} \|f\|_{L^p}$ so $\|\Lambda\| \leq \|\tilde{h}\|_{L^q}$. Also, as before,

$$\begin{aligned} \|\tilde{h}\|_{L^q(\mu)}^q &= \left| \int_{\Omega} \tilde{h} |\tilde{h}|^{q-2} \tilde{h} d\mu \right| = \left| \Lambda\left(|\tilde{h}|^{q-2} \tilde{h}\right) \right| \leq \|\Lambda\| \left(\int_{\Omega} |\tilde{h}|^{q-2} \tilde{h}^p d\mu \right)^{1/p} \\ &= \|\Lambda\| \left(\int_{\Omega} \left(|\tilde{h}|^{q/p}\right)^p \right)^{1/p} = \|\Lambda\| \|h\|_{L^q(\mu)}^{q/p} \end{aligned}$$

and so $\|\tilde{h}\|_{L^q(\mu)} \leq \|\Lambda\|$. It works the same for $p = 1$. Thus $\|\tilde{h}\|_{L^q(\mu)} = \|\Lambda\|$. ■

A situation in which the conditions of the lemma are satisfied is the case where the measure space is σ finite. In fact, you should show this is the only case in which the conditions of the above lemma hold.

Theorem 23.3.6 (Riesz representation theorem) *Let $(\Omega, \mathcal{S}, \mu)$ be σ finite and let $\Lambda \in (L^p(\Omega, \mu))'$, $p \geq 1$. Then there exists a unique $h \in L^q(\Omega, \mu)$, $L^\infty(\Omega, \mu)$ if $p = 1$ such that $\Lambda f = \int h f d\mu$. Also $\|h\| = \|\Lambda\|$. ($\|h\| = \|h\|_q$ if $p > 1$, $\|h\|_\infty$ if $p = 1$). Here $\frac{1}{p} + \frac{1}{q} = 1$.*

Proof: Without loss of generality, assume $\mu(\Omega) = \infty$. By Proposition 10.13.1, either μ is a finite measure or $\mu(\Omega) = \infty$. These are the only two cases. Then let $\{\Omega_n\}$ be a sequence

of disjoint elements of \mathcal{S} having the property that $1 < \mu(\Omega_n) < \infty$, $\cup_{n=1}^{\infty} \Omega_n = \Omega$. Define $r(x) = \sum_{n=1}^{\infty} \frac{1}{n^2} \chi_{\Omega_n}(x) \mu(\Omega_n)^{-1}$, $\tilde{\mu}(E) = \int_E r d\mu$. Thus $\int_{\Omega} r d\mu = \tilde{\mu}(\Omega) = \sum_{n=1}^{\infty} \frac{1}{n^2} < \infty$ so $\tilde{\mu}$ is a finite measure. The above lemma gives the existence part of the conclusion of the theorem. Uniqueness is done as before. ■

With the Riesz representation theorem, it is easy to show that $L^p(\Omega)$, $p > 1$ is a reflexive Banach space. Recall Definition 21.2.14 on Page 546 for the definition. From this, one obtains a weak compactness result.

23.4 Weak Compactness

Theorem 23.4.1 *For $(\Omega, \mathcal{S}, \mu)$ a σ finite measure space and $p > 1$, $L^p(\Omega)$ is reflexive. Thus every bounded sequence has a weakly convergent subsequence.*

Proof: Let $\delta_r : (L^r(\Omega))' \rightarrow L^r(\Omega)$ be defined for $\frac{1}{r} + \frac{1}{r'} = 1$ by $\int (\delta_r \Lambda) g d\mu = \Lambda g$ for all $g \in L^r(\Omega)$. From Theorem 23.3.6 δ_r is one to one, onto, continuous and linear. By the open map theorem, δ_r^{-1} is also one to one, onto, and continuous ($\delta_r \Lambda$ equals the representer of Λ). Thus δ_r^* is also one to one, onto, and continuous by Corollary 21.2.11. Now observe that $J = \delta_p^* \circ \delta_q^{-1}$. To see this, let $z^* \in (L^q)'$, $y^* \in (L^p)'$,

$$\delta_p^* \circ \delta_q^{-1}(\delta_q z^*)(y^*) = (\delta_p^* z^*)(y^*) = z^*(\delta_p y^*) = \int (\delta_q z^*)(\delta_p y^*) d\mu,$$

$$J(\delta_q z^*)(y^*) = y^*(\delta_q z^*) = \int (\delta_p y^*)(\delta_q z^*) d\mu.$$

Therefore $\delta_p^* \circ \delta_q^{-1} = J$ on $\delta_q(L^q)' = L^p$. But the two δ maps are onto and so J is also onto. ■

What about weak compactness in $L^1(\Omega)$? I will give a simple sufficient condition in the case of a finite measure space. More can be said. See for example Dunford and Schwartz [16]. I have this in my Topics in Analysis book also. Recall Proposition 10.9.6 on Page 293 which says equi-integrable is the same as bounded and uniformly integrable. Thus in the following, you can replace equi-integrable with bounded and uniformly integrable.

Theorem 23.4.2 *Let $(\Omega, \mathcal{F}, \mu)$ be a finite measure space and let $\{f_n\}$ be a sequence in $L^1(\Omega)$ which is equi-integrable. Then there exists a subsequence which converges weakly in $L^1(\Omega)$ to some function f .*

Proof: Let

$$\mathcal{E}_n \equiv \{f_n^{-1}(B(z, r)) : r \text{ is a positive rational and } z \in \mathbb{Q} + i\mathbb{Q}\}.$$

Let $\mathcal{E} = \cup_{n=1}^{\infty} \mathcal{E}_n$. Thus \mathcal{E} and \mathcal{E}_n are countable. Also, every open set is the countable union of these sets $B(z, r)$. Now let \mathcal{K} be all finite intersections of sets of \mathcal{E} and include \emptyset and Ω in \mathcal{K} . Then $\sigma(\mathcal{K})$ contains inverse images of Borel sets for each f_n . Thus each f_n is measurable with respect to $\sigma(\mathcal{K})$. Also \mathcal{K} is countable. Then, using a Cantor diagonal argument, we can have $\int_E g_n d\mu$ converges for all $E \in \mathcal{K}$. Let \mathcal{G} be those sets $G \in \sigma(\mathcal{K})$ such that $\int_G g_n d\mu$ converges. Suppose G_k are disjoint and each in \mathcal{G} . Then, since Ω has finite measure, $\lim_{n \rightarrow \infty} \mu(\cup_{j=k}^{\infty} G_j) = 0$ because $\sum_k \mu(G_k)$ converges. Let $G = \cup_{k=1}^{\infty} G_k$ and so $\int_G g_n d\mu - \int_G g_m d\mu = \int_G \sum_{k=1}^{\infty} \chi_{G_k} (g_n - g_m) d\mu = \sum_{k=1}^N \int_{G_k} (g_n - g_m) d\mu +$

$e(N, n, m)$ where $|e(N, n, m)| < \varepsilon$ for all n, m provided N is chosen large enough. This is by uniform integrability which is a consequence of equi-integrability. See Proposition 10.9.6. It follows that

$$\begin{aligned} \left| \int_G g_n d\mu - \int_G g_m d\mu \right| &\leq \left| \sum_{k=1}^N \int_{G_k} (g_n - g_m) d\mu \right| + |e(N, n, m)| \\ &< \left| \sum_{k=1}^N \int_{G_k} (g_n - g_m) d\mu \right| + \varepsilon < 2\varepsilon \end{aligned}$$

provided n, m are large enough. Thus \mathcal{G} is closed with respect to countable disjoint unions. If $\int_G g_n d\mu$ converges, then $\int_{G^c} g_n d\mu = \int_\Omega g_n d\mu - \int_G g_n d\mu$ and so $\int_{G^c} g_n d\mu$ converges. Hence, by Dynkin's lemma, $\mathcal{G} \supseteq \sigma(\mathcal{K})$. For $E \in \sigma(\mathcal{K})$ define

$$\begin{aligned} \lambda(E) &\equiv \lim_{n \rightarrow \infty} \int_E g_n d\mu, \text{ then } \lambda \ll \mu \text{ so there is } g \text{ such that} \\ \int_E g d\mu &= \lambda(E) = \lim_{n \rightarrow \infty} \int_E g_n d\mu \text{ by Radon Nikodym, } g \in L^1 \end{aligned}$$

That λ is a measure follows from the above argument that \mathcal{G} is closed with respect to countable disjoint unions.

Now it was just shown that for s a simple function measurable with respect to $\sigma(\mathcal{K})$,

$$\int s g d\mu = \lim_{n \rightarrow \infty} \int s g_n d\mu.$$

Can we replace s with $h \in L^\infty(\Omega, \sigma(\mathcal{K}), \mu)$? Letting h be a representative which is uniformly bounded, there exists a sequence of simple functions $\{s_n\}$ which converges uniformly to h .

$$\begin{aligned} \left| \int h g d\mu - \int h g_n d\mu \right| &\leq \left| \int h g d\mu - \int s g d\mu \right| \\ &+ \left| \int s g d\mu - \int s g_n d\mu \right| + \left| \int s g_n d\mu - \int h g_n d\mu \right| \end{aligned}$$

The first term on the right is no more than $\varepsilon \|g\|_{L^1}$ because s was chosen to be uniformly within ε of h . As to the last term, it is no more than $\varepsilon \max_n \|g_n\|_{L^1}$ which is no more than εC since the equi-integrability implies $\|g_n\|_{L^1}$ is bounded. The middle term converges to 0 and so $\lim_{n \rightarrow \infty} |\int h g d\mu - \int h g_n d\mu| = 0$.

Now consider $L^\infty(\Omega, \sigma(\mathcal{K}), \mu) \xleftarrow{i^*} L^\infty(\Omega, \mathcal{F}, \mu)$ where the inclusion map i is $L^1(\Omega, \sigma(\mathcal{K}), \mu) \xrightarrow{i} L^1(\Omega, \mathcal{F}, \mu)$ continuous. Let $h \in L^\infty(\Omega, \mathcal{F}, \mu)$ so $i^* h \in L^\infty(\Omega, \sigma(\mathcal{K}), \mu)$. Then

$$\lim_{n \rightarrow \infty} \int h g_n d\mu = \lim_{n \rightarrow \infty} \int i g_n d\mu = \lim_{n \rightarrow \infty} \int i^* h g_n d\mu = \lim_{n \rightarrow \infty} \int i^* h g d\mu = \lim_{n \rightarrow \infty} \int h g d\mu$$

and this shows that g_n converges weakly to g . ■

One can extend this to an arbitrary measure space by fussing with more details that involve consideration of a σ algebra which is σ finite.

23.5 The Dual Space of $L^\infty(\Omega)$

What about the dual space of $L^\infty(\Omega)$? This will involve the following Lemma. Also recall the notion of total variation defined in Definition 23.2.2.

Lemma 23.5.1 *Let (Ω, \mathcal{F}) be a measure space. Denote by $BV(\Omega)$ the space of **finitely** additive complex measures ν such that $|\nu|(\Omega) < \infty$. This means that if $\{E_i\}_{i=1}^n$ is disjoint, then $\nu(\cup_{i=1}^n E_i) = \sum_{i=1}^n \nu(E_i)$ for any $n \in \mathbb{N}$. Then defining $\|\nu\| \equiv |\nu|(\Omega)$, it follows that $BV(\Omega)$ is a Banach space.*

Proof: It is obvious that $BV(\Omega)$ is a vector space with the obvious conventions involving scalar multiplication. Why is $\|\cdot\|$ a norm? All the axioms are obvious except for the triangle inequality. However, this is not too hard either.

$$\begin{aligned} \|\mu + \nu\| &\equiv |\mu + \nu|(\Omega) = \sup_{\pi(\Omega)} \left\{ \sum_{A \in \pi(\Omega)} |\mu(A) + \nu(A)| \right\} \\ &\leq \sup_{\pi(\Omega)} \left\{ \sum_{A \in \pi(\Omega)} |\mu(A)| \right\} + \sup_{\pi(\Omega)} \left\{ \sum_{A \in \pi(\Omega)} |\nu(A)| \right\} \equiv |\mu|(\Omega) + |\nu|(\Omega) = \|\nu\| + \|\mu\|. \end{aligned}$$

Suppose now that $\{\nu_n\}$ is a Cauchy sequence. For each $E \in \mathcal{F}$, $|\nu_n(E) - \nu_m(E)| \leq \|\nu_n - \nu_m\|$ and so the sequence of complex numbers $\nu_n(E)$ converges. That to which it converges is called $\nu(E)$. Then it is obvious that $\nu(E)$ is finitely additive. Why is $|\nu|$ finite? Since $\|\cdot\|$ is a norm, it follows that there exists a constant C such that for all n , $|\nu_n|(\Omega) < C$. Let $\pi(\Omega)$ be any partition. Then $\sum_{A \in \pi(\Omega)} |\nu(A)| = \lim_{n \rightarrow \infty} \sum_{A \in \pi(\Omega)} |\nu_n(A)| \leq C$. Hence $\nu \in BV(\Omega)$. Let $\varepsilon > 0$ be given and let N be such that if $n, m > N$, then $\|\nu_n - \nu_m\| < \varepsilon/2$. Pick any such n . Then choose $\pi(\Omega)$ such that

$$\begin{aligned} |\nu - \nu_n|(\Omega) - \varepsilon/2 &< \sum_{A \in \pi(\Omega)} |\nu(A) - \nu_n(A)| \\ &= \lim_{m \rightarrow \infty} \sum_{A \in \pi(\Omega)} |\nu_m(A) - \nu_n(A)| < \lim_{m \rightarrow \infty} \inf_{m \rightarrow \infty} |\nu_n - \nu_m|(\Omega) \leq \varepsilon/2 \end{aligned}$$

It follows that $\lim_{n \rightarrow \infty} \|\nu - \nu_n\| = 0$. ■

Corollary 23.5.2 *Suppose (Ω, \mathcal{F}) is a measure space as above and suppose μ is a measure defined on \mathcal{F} . Denote by $BV(\Omega; \mu)$ those finitely additive measures of $BV(\Omega)$ ν such that $\nu \ll \mu$ in the usual sense that if $\mu(E) = 0$, then $\nu(E) = 0$. Then $BV(\Omega; \mu)$ is a closed subspace of $BV(\Omega)$.*

Proof: It is clear that it is a subspace. Is it closed? Suppose $\nu_n \rightarrow \nu$ and each ν_n is in $BV(\Omega; \mu)$. Then if $\mu(E) = 0$, it follows that $\nu_n(E) = 0$ and so $\nu(E) = 0$ also, being the limit of 0. ■

Definition 23.5.3 *For s a simple function $s(\omega) = \sum_{k=1}^n c_k \mathcal{X}_{E_k}(\omega)$ and $\nu \in BV(\Omega)$, define an “integral” with respect to ν as follows. $\int s d\nu \equiv \sum_{k=1}^n c_k \nu(E_k)$. For f function which is in $L^\infty(\Omega; \mu)$, define $\int f d\nu$ as follows. Applying Theorem 9.1.6, to the positive and negative parts of real and imaginary parts of f , there exists a sequence of simple functions $\{s_n\}$ which converges uniformly to f off a set of μ measure zero. Then $\int f d\nu \equiv \lim_{n \rightarrow \infty} \int s_n d\nu$*

Lemma 23.5.4 *The above definition of the integral with respect to a finitely additive measure in $BV(\Omega; \mu)$ is well defined.*

Proof: First consider the claim about the integral being well defined on the simple functions. This is clearly true if it is required that the c_k are disjoint and the E_k also disjoint having union equal to Ω . Thus define the integral of a simple function in this manner. First write the simple function as $\sum_{k=1}^n c_k \mathcal{X}_{E_k}$ where the c_k are the values of the simple function. Then use the above formula to define the integral. Next suppose the E_k are disjoint but the c_k are not necessarily distinct. Let the distinct values of the c_k be a_1, \dots, a_m

$$\begin{aligned} \sum_k c_k \mathcal{X}_{E_k} &= \sum_j a_j \left(\sum_{i: c_i = a_j} \mathcal{X}_{E_i} \right) = \sum_j a_j \mathbf{v}(\cup_{i: c_i = a_j} E_i) \\ &= \sum_j a_j \sum_{i: c_i = a_j} \mathbf{v}(E_i) = \sum_k c_k \mathbf{v}(E_k) \end{aligned}$$

and so the same formula for the integral of a simple function is obtained in this case also. Now consider two simple functions $s = \sum_{k=1}^n a_k \mathcal{X}_{E_k}$, $t = \sum_{j=1}^m b_j \mathcal{X}_{F_j}$ where the a_k and b_j are the distinct values of the simple functions. Then from what was just shown,

$$\begin{aligned} \int (\alpha s + \beta t) d\mathbf{v} &= \int \left(\sum_{k=1}^n \sum_{j=1}^m \alpha a_k \mathcal{X}_{E_k \cap F_j} + \sum_{j=1}^m \sum_{k=1}^n \beta b_j \mathcal{X}_{E_k \cap F_j} \right) d\mathbf{v} \\ &= \int \left(\sum_{j,k} \alpha a_k \mathcal{X}_{E_k \cap F_j} + \beta b_j \mathcal{X}_{E_k \cap F_j} \right) d\mathbf{v} = \sum_{j,k} (\alpha a_k + \beta b_j) \mathbf{v}(E_k \cap F_j) \\ &= \sum_{k=1}^n \sum_{j=1}^m \alpha a_k \mathbf{v}(E_k \cap F_j) + \sum_{j=1}^m \sum_{k=1}^n \beta b_j \mathbf{v}(E_k \cap F_j) \\ &= \sum_{k=1}^n \alpha a_k \mathbf{v}(E_k) + \sum_{j=1}^m \beta b_j \mathbf{v}(F_j) = \alpha \int s d\mathbf{v} + \beta \int t d\mathbf{v} \end{aligned}$$

Thus the integral is linear on simple functions so, in particular, the formula given in the above definition is well defined regardless.

So what about the definition for $f \in L^\infty(\Omega; \mu)$? Since $f \in L^\infty$, there is a set of μ measure zero N such that on N^C there exists a sequence of simple functions which converges uniformly to f on N^C . Consider s_n and s_m . As in the above, they can be written as $\sum_{k=1}^p c_k^n \mathcal{X}_{E_k}$, $\sum_{k=1}^p c_k^m \mathcal{X}_{E_k}$ respectively, where the E_k are disjoint having union equal to Ω . Then by uniform convergence, if m, n are sufficiently large, $|c_k^n - c_k^m| < \varepsilon$ or else the corresponding E_k is contained in N^C a set of \mathbf{v} measure 0 thanks to $\mathbf{v} \ll \mu$. Hence

$$\begin{aligned} \left| \int s_n d\mathbf{v} - \int s_m d\mathbf{v} \right| &= \left| \sum_{k=1}^p (c_k^n - c_k^m) \mathbf{v}(E_k) \right| \\ &\leq \sum_{k=1}^p |c_k^n - c_k^m| \mathbf{v}(E_k) \leq \varepsilon \|\mathbf{v}\| \end{aligned}$$

and so the integrals of these simple functions converge. Similar reasoning shows that the definition is not dependent on the choice of approximating sequence. ■

Note also that for s simple, $|\int s d\nu| \leq \|s\|_{L^\infty} \|\nu\|(\Omega) = \|s\|_{L^\infty} \|\nu\|$

Next the dual space of $L^\infty(\Omega; \mu)$ will be identified with $BV(\Omega; \mu)$. First here is a simple observation. Let $\nu \in BV(\Omega; \mu)$. Then define the following for $f \in L^\infty(\Omega; \mu)$.

$$T_\nu(f) \equiv \int f d\nu$$

Lemma 23.5.5 For T_ν just defined, $|T_\nu f| \leq \|f\|_{L^\infty} \|\nu\|$.

Proof: As noted above, the conclusion true if f is simple. Now if f is in L^∞ , then it is the uniform limit of simple functions off a set of μ measure zero. Therefore, by the definition of the T_ν ,

$$|T_\nu f| = \lim_{n \rightarrow \infty} |T_\nu s_n| \leq \lim_{n \rightarrow \infty} \|s_n\|_{L^\infty} \|\nu\| = \|f\|_{L^\infty} \|\nu\|. \blacksquare$$

Thus each T_ν is in $(L^\infty(\Omega; \mu))'$. \blacksquare

Here is the representation theorem, due to Kantorovitch, which describes the dual of $L^\infty(\Omega; \mu)$.

Theorem 23.5.6 Let $\theta : BV(\Omega; \mu) \rightarrow (L^\infty(\Omega; \mu))'$ be given by $\theta(\nu) \equiv T_\nu$. Then θ is one to one, onto and preserves norms.

Proof: It was shown in the above lemma that θ maps into $(L^\infty(\Omega; \mu))'$. It is obvious that θ is linear. Why does it preserve norms? From the above lemma,

$$\|\theta \nu\| \equiv \sup_{\|f\|_\infty \leq 1} |T_\nu f| \leq \|\nu\|$$

It remains to turn the inequality around. Let $\pi(\Omega)$ be a partition. Then

$$\sum_{A \in \pi(\Omega)} |\nu(A)| = \sum_{A \in \pi(\Omega)} \text{sgn}(\nu(A)) \nu(A) \equiv \int f d\nu$$

where $\text{sgn}(\nu(A))$ is defined to be a complex number of modulus 1, $\text{sgn}(\nu(A)) \nu(A) = |\nu(A)|$ and

$$f(\omega) = \sum_{A \in \pi(\Omega)} \text{sgn}(\nu(A)) \mathcal{X}_A(\omega).$$

Therefore, choosing $\pi(\Omega)$ suitably, since $\|f\|_\infty \leq 1$,

$$\begin{aligned} \|\nu\| - \varepsilon &= |\nu|(\Omega) - \varepsilon \leq \sum_{A \in \pi(\Omega)} |\nu(A)| = T_\nu(f) \\ &= |T_\nu(f)| = |\theta(\nu)(f)| \leq \|\theta(\nu)\| \leq \|\nu\| \end{aligned}$$

Thus θ preserves norms. Hence it is one to one also. Why is θ onto?

Let $\Lambda \in (L^\infty(\Omega; \mu))'$. Then define

$$\nu(E) \equiv \Lambda(\mathcal{X}_E) \tag{23.7}$$

This is obviously finitely additive because Λ is linear. Also, if $\mu(E) = 0$, then $\mathcal{X}_E = 0$ in L^∞ and so $\Lambda(\mathcal{X}_E) = 0$. If $\pi(\Omega)$ is any partition of Ω , then

$$\begin{aligned} \sum_{A \in \pi(\Omega)} |\nu(A)| &= \sum_{A \in \pi(\Omega)} |\Lambda(\mathcal{X}_A)| = \sum_{A \in \pi(\Omega)} \text{sgn}(\Lambda(\mathcal{X}_A)) \Lambda(\mathcal{X}_A) \\ &= \Lambda\left(\sum_{A \in \pi(\Omega)} \text{sgn}(\Lambda(\mathcal{X}_A)) \mathcal{X}_A\right) \leq \|\Lambda\| \end{aligned}$$

and so $\|v\| \leq \|\Lambda\|$ showing that $v \in BV(\Omega; \mu)$. Also from 23.7, if $s = \sum_{k=1}^n c_k \mathcal{X}_{E_k}$ is a simple function,

$$\int s dv = \sum_{k=1}^n c_k v(E_k) = \sum_{k=1}^n c_k \Lambda(\mathcal{X}_{E_k}) = \Lambda\left(\sum_{k=1}^n c_k \mathcal{X}_{E_k}\right) = \Lambda(s)$$

Then letting $f \in L^\infty(\Omega; \mu)$, there exists a sequence of simple functions converging to f uniformly off a set of μ measure zero and so passing to a limit in the above with s replaced with s_n it follows that $\Lambda(f) = \int f dv$ and so θ is onto. ■

23.6 Non σ Finite Case

It turns out that for $p > 1$, you don't have to assume the measure space is σ finite. The Riesz representation theorem holds always. The proof involves the notion of uniform convexity. First recall Clarkson's inequalities. These fundamental inequalities were used to verify that $L^p(\Omega)$ is uniformly convex. More precisely, the unit ball in $L^p(\Omega)$ is uniformly convex.

Lemma 23.6.1 *Let $2 \leq p$. Then*

$$\left\| \frac{f+g}{2} \right\|_{L^p}^p + \left\| \frac{f-g}{2} \right\|_{L^p}^p \leq \frac{1}{2} (\|f\|_{L^p}^p + \|g\|_{L^p}^p)$$

Let $1 < p < 2$. then for $1/p + 1/q = 1$,

$$\left\| \frac{f+g}{2} \right\|_{L^p}^q + \left\| \frac{f-g}{2} \right\|_{L^p}^q \leq \left(\frac{1}{2} \|f\|_{L^p}^p + \frac{1}{2} \|g\|_{L^p}^p \right)^{q/p}$$

Recall the following definition of uniform convexity.

Definition 23.6.2 *A Banach space, X , is said to be uniformly convex if whenever $\|x_n\| \leq 1$ and $\left\| \frac{x_n + x_m}{2} \right\| \rightarrow 1$ as $n, m \rightarrow \infty$, then $\{x_n\}$ is a Cauchy sequence and $x_n \rightarrow x$ where $\|x\| = 1$. More precisely, for every $\varepsilon > 0$, there is a $\delta > 0$ such that if $\|x + y\| > 2 - \delta$ for $\|x\|, \|y\|$ both 1, then $\|x - y\| < \varepsilon$.*

Observe that Clarkson's inequalities imply L^p is uniformly convex for all $p > 1$. Consider the harder case where $1 < p$. The other case is similar. Say $\|f\| = \|g\| = 1$ and $\|f + g\|_{L^p} > 2 - \delta$. Then from the second inequality $\left(\frac{2-\delta}{2} \right)^q + \left\| \frac{f-g}{2} \right\|_{L^p}^q \leq 1$ and so

$$\|f - g\|_{L^p}^q \leq 2^q \left(1 - \left(\frac{2-\delta}{2} \right)^q \right) < \varepsilon$$

provided δ is small enough.

Uniformly convex spaces have a very nice property which is described in the following lemma. Roughly, this property is that any element of the dual space achieves its norm at some point of the closed unit ball.

Lemma 23.6.3 *Let X be uniformly convex and let $\phi \in X'$. Then there exists $x \in X$ such that $\|x\| = 1$, $\phi(x) = \|\phi\|$.*

Proof: There is nothing to show if $\phi = 0$ so suppose it is not. Let $\|x_n\| = 1$ and let $\phi(x_n) \rightarrow \|\phi\|$. Then as $n, m \rightarrow \infty$, $\phi\left(\frac{x_n + x_m}{2}\right) \rightarrow \|\phi\|$. Without loss of generality, we can also assume $\phi(x_n)$ is positive. Hence if m, n are large enough, then

$$\|\phi\|(1 - \varepsilon) < \phi\left(\frac{x_n + x_m}{2}\right) \leq \|\phi\| \left\| \frac{x_n + x_m}{2} \right\|$$

Thus, if m, n are large enough, $\|x_n + x_m\| \geq 2(1 - \varepsilon)$. It follows that $\lim_{m, n \rightarrow \infty} \|x_n + x_m\| = 2$ and so by uniform convexity, $\lim_{m, n \rightarrow \infty} \|x_n - x_m\| = 0$. Thus the sequence is a Cauchy sequence and so there is x , $\|x\| = 1$ and $x_n \rightarrow x$ so $\|\phi\| = \lim_{n \rightarrow \infty} \phi(x_n) = \phi(x)$. ■

The proof of the Riesz representation theorem will be based on the following lemma which says that if you can show a directional derivative exists, then it can be used to represent a functional in terms of this directional derivative. It is very interesting for its own sake.

Lemma 23.6.4 (McShane) *Let X be a complex normed linear space and let $\phi \in X'$. Suppose there exists $x \in X$, $\|x\| = 1$ with $\phi(x) = \|\phi\| \neq 0$. Let $y \in X$ and let $\psi_y(t) = \|x + ty\|$ for $t \in \mathbb{R}$. Suppose $\psi'_y(0)$ exists for each $y \in X$. Then for all $y \in X$,*

$$\psi'_y(0) + i\psi'_{-iy}(0) = \|\phi\|^{-1} \phi(y).$$

Proof: Suppose first that $\|\phi\| = 1$. The idea is to show that in the limit as $t \rightarrow 0$,

$$\frac{|1 + t\phi(y)| - 1}{t}, \frac{\|x + ty\| - \|x\|}{t}$$

act the same. The first part of the argument is devoted to showing this.

By assumption, there is x such that $\|x\| = 1$ and $\phi(x) = 1 = \|\phi\|$. Then $\phi(y - \phi(y)x) = 0$ and so

$$\phi(x + t(y - \phi(y)x)) = \phi(x) + t\phi(y) - t\phi(y)\phi(x) \stackrel{=1}{=} \phi(x) = 1 = \|\phi\|.$$

Therefore, $\|x + t(y - \phi(y)x)\| \geq 1$ since, from the above,

$$\|\phi\| \|x + t(y - \phi(y)x)\| = \|x + t(y - \phi(y)x)\| \geq \|\phi\| = 1$$

Also for small t , $|\phi(y)t| < 1$, and so $1 \leq \|x + t(y - \phi(y)x)\| = \|(1 - \phi(y)t)x + ty\|$

$$\leq |1 - \phi(y)t| \left\| x + \frac{t}{1 - \phi(y)t} y \right\|.$$

Divide both sides by $|1 - \phi(y)t|$. Using the standard formula for the sum of a geometric series,

$$1 + t\phi(y) + o(t) = \frac{1}{1 - t\phi(y)}$$

Therefore,

$$\frac{1}{|1 - \phi(y)t|} = |1 + \phi(y)t + o(t)| \leq \left\| x + \frac{t}{1 - \phi(y)t} y \right\| = \|x + ty + o(t)\| \quad (23.8)$$

where $\lim_{t \rightarrow 0} o(t)(t^{-1}) = 0$. Thus, $|1 + \phi(y)t| \leq \|x + ty\| + o(t)$. Now since $t\phi(y) \in \mathbb{C}$,

$$|1 + t\phi(y)| - 1 \geq 1 + t\operatorname{Re} \phi(y) - 1 = t\operatorname{Re} \phi(y).$$

Thus for $t > 0$,

$$\operatorname{Re} \phi(y) \leq \frac{|1 + t\phi(y)| - 1}{t} \stackrel{\|x\|=1}{\leq} \frac{\|x + ty\| - \|x\|}{t} + \frac{o(t)}{t}$$

and for $t < 0$,

$$\operatorname{Re} \phi(y) \geq \frac{|1 + t\phi(y)| - 1}{t} \geq \frac{\|x + ty\| - \|x\|}{t} + \frac{o(t)}{t}$$

By assumption that the directional derivative exists, and letting $t \rightarrow 0+$ and $t \rightarrow 0-$,

$$\operatorname{Re} \phi(y) = \lim_{t \rightarrow 0} \frac{\|x + ty\| - \|x\|}{t} = \psi'_y(0).$$

Now $\phi(y) = \operatorname{Re} \phi(y) + i \operatorname{Im} \phi(y)$ so $\phi(-iy) = -i(\phi(y)) = -i \operatorname{Re} \phi(y) + \operatorname{Im} \phi(y)$ and

$$\phi(-iy) = \operatorname{Re} \phi(-iy) + i \operatorname{Im} \phi(-iy).$$

Hence $\operatorname{Re} \phi(-iy) = \operatorname{Im} \phi(y)$. Consequently,

$$\begin{aligned} \phi(y) &= \operatorname{Re} \phi(y) + i \operatorname{Im} \phi(y) = \operatorname{Re} \phi(y) + i \operatorname{Re} \phi(-iy) \\ &= \psi'_y(0) + i \psi'_{-iy}(0). \end{aligned}$$

This proves the lemma when $\|\phi\| = 1$. For arbitrary $\phi \neq 0$, let $\phi(x) = \|\phi\|$, $\|x\| = 1$. Then from above, if $\phi_1(y) \equiv \|\phi\|^{-1} \phi(y)$, $\|\phi_1\| = 1$ and so from what was just shown,

$$\phi_1(y) = \frac{\phi(y)}{\|\phi\|} = \psi'_y(0) + i \psi'_{-iy}(0) \blacksquare$$

This shows you can represent ϕ in terms of this directional derivative.

Now here are some short observations. For $t \in \mathbb{R}$, $p > 1$, and $x, y \in \mathbb{C}$, $x \neq 0$

$$\begin{aligned} \lim_{t \rightarrow 0} \frac{|x + ty|^p - |x|^p}{t} &= p|x|^{p-2} (\operatorname{Re} x \operatorname{Re} y + \operatorname{Im} x \operatorname{Im} y) \\ &= p|x|^{p-2} \operatorname{Re}(\bar{x}y) \end{aligned} \tag{23.9}$$

Also from convexity of $f(r) = r^p$, for $|t| < 1$,

$$\begin{aligned} |x + ty|^p - |x|^p &\leq \|x\| + |t| \|y\|^p - |x|^p \\ &= \left[(1 + |t|) \left(\frac{|x| + |t||y|}{1 + |t|} \right) \right]^p - |x|^p \leq (1 + |t|)^p \frac{|x|^p}{1 + |t|} + \frac{|t||y|^p}{1 + |t|} - |x|^p \\ &\leq (1 + |t|)^{p-1} (|x|^p + |t||y|^p) - |x|^p \leq \left((1 + |t|)^{p-1} - 1 \right) |x|^p + 2^{p-1} |t||y|^p \end{aligned}$$

Now for $f(t) \equiv (1 + t)^{p-1}$, $f'(t)$ is uniformly bounded, depending on p , for $t \in [0, 1]$. Hence the above is dominated by an expression of the form

$$C_p (|x|^p + |y|^p) |t| \tag{23.10}$$

The above lemma and uniform convexity of L^p can be used to prove a general version of the Riesz representation theorem next. Let $p > 1$ and let $\eta : L^q \rightarrow (L^p)'$ be defined by

$$\eta(g)(f) = \int_{\Omega} gf \, d\mu. \tag{23.11}$$

Theorem 23.6.5 (*Riesz representation theorem $p > 1$*) The map η is 1-1, onto, continuous, and

$$\|\eta g\| = \|g\|, \quad \|\eta\| = 1.$$

Proof: Obviously η is linear. Suppose $\eta g = 0$. Then $0 = \int g f d\mu$ for all $f \in L^p$. Let $f = |g|^{q-2}\bar{g}$. Then $f \in L^p$ and so $0 = \int |g|^q d\mu$. Hence $g = 0$ and η is one to one. That $\eta g \in (L^p)'$ is obvious from the Holder inequality. In fact, $|\eta(g)(f)| \leq \|g\|_q \|f\|_p$, and so $\|\eta(g)\| \leq \|g\|_q$. To see that equality holds, let $f = |g|^{q-2}\bar{g} \|g\|_q^{1-q}$. Then $\|f\|_p = 1$ and

$$\eta(g)(f) = \int_{\Omega} |g|^q d\mu \|g\|_q^{1-q} = \|g\|_q.$$

Thus $\|\eta\| = 1$.

It remains to show η is onto. Let $\phi \in (L^p)'$. Is $\phi = \eta g$ for some $g \in L^q$? Without loss of generality, assume $\phi \neq 0$. By uniform convexity of L^p , Lemma 23.6.3, there exists g such that $\phi g = \|\phi\|$, $g \in L^p$, $\|g\| = 1$. For $f \in L^p$, define $\phi_f(t) \equiv \int_{\Omega} |g + tf|^p d\mu$. Thus $\psi_f(t) \equiv \|g + tf\|_p \equiv \phi_f(t)^{\frac{1}{p}}$. Does $\phi'_f(0)$ exist? Let $[g = 0]$ denote the set $\{x : g(x) = 0\}$.

$$\frac{\phi_f(t) - \phi_f(0)}{t} = \int \frac{(|g + tf|^p - |g|^p)}{t} d\mu$$

From 23.10, the integrand is bounded by $C_p(|f|^p + |g|^p)$. Therefore, using 23.9, the dominated convergence theorem applies and it follows $\phi'_f(0) =$

$$\begin{aligned} \lim_{t \rightarrow 0} \frac{\phi_f(t) - \phi_f(0)}{t} &= \lim_{t \rightarrow 0} \left[\int_{[g=0]} |t|^{p-1} |f|^p d\mu + \int_{[g \neq 0]} \frac{(|g + tf|^p - |g|^p)}{t} d\mu \right] \\ &= p \int_{[g \neq 0]} |g|^{p-2} \operatorname{Re}(\bar{g}f) d\mu = p \int |g|^{p-2} \operatorname{Re}(\bar{g}f) d\mu \end{aligned}$$

Hence $\psi'_f(0) = \|g\|^{-\frac{p}{q}} \int |g(x)|^{p-2} \operatorname{Re}(g(x)\bar{f}(x)) d\mu$. Note $\frac{1}{p} - 1 = -\frac{1}{q}$. Therefore,

$$\psi'_{-if}(0) = \|g\|^{-\frac{p}{q}} \int |g(x)|^{p-2} \operatorname{Re}(ig(x)\bar{f}(x)) d\mu.$$

But $\operatorname{Re}(ig\bar{f}) = \operatorname{Im}(-g\bar{f})$ and so by the McShane lemma,

$$\begin{aligned} \phi(f) &= \|\phi\| \|g\|^{-\frac{p}{q}} \int |g(x)|^{p-2} [\operatorname{Re}(g(x)\bar{f}(x)) + i \operatorname{Re}(ig(x)\bar{f}(x))] d\mu \\ &= \|\phi\| \|g\|^{-\frac{p}{q}} \int |g(x)|^{p-2} [\operatorname{Re}(g(x)\bar{f}(x)) + i \operatorname{Im}(-g(x)\bar{f}(x))] d\mu \\ &= \|\phi\| \|g\|^{-\frac{p}{q}} \int |g(x)|^{p-2} \bar{g}(x)f(x) d\mu. \end{aligned}$$

This shows that $\phi = \eta(\|\phi\| \|g\|^{-\frac{p}{q}} |g|^{p-2}\bar{g})$ and verifies η is onto. ■

23.7 The Dual Space of $C_0(X)$

Consider the dual space of $C_0(X)$ where X is a locally compact Hausdorff space. It will turn out to be a space of measures. To show this, the following lemma will be convenient. Recall this space is defined as follows.

Definition 23.7.1 $f \in C_0(X)$ means that for every $\varepsilon > 0$ there exists a compact set K such that $|f(x)| < \varepsilon$ whenever $x \notin K$. Recall the norm on this space is

$$\|f\|_\infty \equiv \|f\| \equiv \sup \{|f(x)| : x \in X\}$$

Also define $C_0^+(X)$ to be the nonnegative functions of $C_0(X)$.

From the representation theorem about positive linear functionals on $C_0(X)$, we know that if Λ is such a positive linear functional, then $\Lambda f = \int_X f d\mu$. What if Λ is also continuous so that $|\Lambda f| \leq \|\Lambda\| \|f\|_\infty$?

Lemma 23.7.2 Suppose $\Lambda : C_0(X) \rightarrow \mathbb{C}$ is a positive linear functional which is also continuous. Then if μ is the Radon measure representing Λ , it follows that $\|\Lambda\| = \mu(X)$ so in particular, μ is finite.

Proof: From the regularity of μ , $\mu(X) = \sup \{\mu(K) : K \subseteq X, K \text{ compact}\}$. For such a K let $K \prec f \prec X$ and so it follows that $\mu(X) = \sup \{\Lambda f : f \prec X\}$. However, $0 \leq \Lambda f \leq \|\Lambda\| \|f\|_\infty \leq \|\Lambda\|$ and so $\mu(X) \leq \|\Lambda\|$. To go the other direction, use density of $C_c(X)$ in $C_0(X)$ to obtain $f \in C_c(X)$ such that $\|f\|_\infty \leq 1$ and $\|\Lambda\| < |\Lambda f| + \varepsilon$. Since Λ is a positive linear functional, one can assume that $f \geq 0$ since otherwise the inequality could be strengthened by replacing f with its positive part f^+ . Then with this f ,

$$\mu(X) \leq \|\Lambda\| < |\Lambda f| + \varepsilon = \Lambda f + \varepsilon \leq \mu(X) + \varepsilon$$

and so, since ε is arbitrary, $\mu(X) = \|\Lambda\|$. ■

Next consider the case where L is in $\mathcal{L}(C_0(X), \mathbb{C})$ so it is not known to take nonnegative functions to nonnegative scalars.

Lemma 23.7.3 Let $L \in \mathcal{L}(C_0(X), \mathbb{C})$. Then there exists $\lambda : C_0^+(X) \rightarrow [0, \infty)$ which satisfies

$$\begin{aligned} \lambda(af + bg) &= a\lambda(f) + b\lambda(g), \text{ if } a, b \geq 0 \\ |\lambda(f)| &\leq \|L\| \|f\|_\infty \end{aligned} \quad (23.12)$$

Proof: Define, for $f \in C_0^+(X)$, $\lambda(f) \equiv \sup \{|Lg| : |g| \leq f\}$. Then the second part is obvious because

$$|\lambda(f)| = \lambda(f) \leq \sup \{\|L\| \|g\|_\infty\} \leq \|L\| \|f\|_\infty$$

Consider the first claim of 23.12. It is obvious that $\lambda(0f) = 0\lambda(f)$ from the above. If $c > 0$, why is $\lambda(cf) = c\lambda(f)$? If $|g| \leq cf$, then $\frac{1}{c}|g| \leq f$ and so $\frac{1}{c}|Lg| \leq \lambda(f)$ so $|Lg| \leq c\lambda(f)$. Taking sup for all such g , $\lambda(cf) \leq c\lambda(f)$. Thus also $\lambda(f) = \lambda(\frac{1}{c}cf) \leq \frac{1}{c}\lambda(cf)$ so $c\lambda(f) \leq \lambda(cf)$ showing that $c\lambda(f) = \lambda(cf)$ if $c \geq 0$. It remains to verify that $\lambda(f_1 + f_2) = \lambda(f_1) + \lambda(f_2)$.

For $f_j \in C_0^+(X)$, there exists $g_i \in C_0(X)$ such that $|g_i| \leq f_i$ and

$$\begin{aligned} \lambda(f_1) + \lambda(f_2) &\leq |L(g_1)| + |L(g_2)| + 2\varepsilon = L(\omega_1 g_1) + L(\omega_2 g_2) + 2\varepsilon \\ &= L(\omega_1 g_1 + \omega_2 g_2) + 2\varepsilon = |L(\omega_1 g_1 + \omega_2 g_2)| + 2\varepsilon \end{aligned}$$

where $|g_i| \leq f_i$ and $|\omega_i| = 1$ and $\omega_i L(g_i) = |L(g_i)|$. Now

$$|\omega_1 g_1 + \omega_2 g_2| \leq |g_1| + |g_2| \leq f_1 + f_2$$

and so the above shows

$$\lambda(f_1) + \lambda(f_2) \leq \lambda(f_1 + f_2) + 2\varepsilon.$$

Since ε is arbitrary, $\lambda(f_1) + \lambda(f_2) \leq \lambda(f_1 + f_2)$. It remains to verify the other inequality.

Now let $|g| \leq f_1 + f_2$,

$$|Lg| \geq \lambda(f_1 + f_2) - \varepsilon.$$

Let

$$h_i(x) = \begin{cases} \frac{f_i(x)g(x)}{f_1(x)+f_2(x)} & \text{if } f_1(x) + f_2(x) > 0, \\ 0 & \text{if } f_1(x) + f_2(x) = 0. \end{cases}$$

Then h_i is continuous and $h_1(x) + h_2(x) = g(x)$, $|h_i| \leq f_i$. The function h_i is clearly continuous at points x where $f_1(x) + f_2(x) > 0$. The reason it is continuous at a point where $f_1(x) + f_2(x) = 0$ is that at every point y where $f_1(y) + f_2(y) > 0$, the top description of the function gives

$$|h_i(y)| = \left| \frac{f_i(y)g(y)}{f_1(y) + f_2(y)} \right| \leq |g(y)| \leq f_1(y) + f_2(y)$$

so if $|y - x|$ is small enough, $|h_i(y)|$ is also small. Then it follows from this definition of the h_i that

$$\begin{aligned} -\varepsilon + \lambda(f_1 + f_2) &\leq |Lg| = |Lh_1 + Lh_2| \leq |Lh_1| + |Lh_2| \\ &\leq \lambda(f_1) + \lambda(f_2). \end{aligned}$$

Since $\varepsilon > 0$ is arbitrary, this shows that

$$\lambda(f_1 + f_2) \leq \lambda(f_1) + \lambda(f_2) \leq \lambda(f_1 + f_2) \quad \blacksquare$$

λ cannot be linear because it is not defined on a vector space. However, it wants to be linear. This is the content of the above lemma. Therefore, I will call λ righteous.

23.7.1 Extending Righteous Functionals

The process of extending such a righteous functional to one which is linear is the same process used earlier with the abstract Lebesgue integral. It is just like Theorem 10.7.8 except here the functional is defined on continuous functions which are nonnegative rather than measurable nonnegative functions. The inequality of 23.12 is also preserved.

Lemma 23.7.4 Suppose λ is a mapping which has $\lambda(f) \geq 0$ which is defined on $C_0^+(X)$ such that

$$\lambda(af + bg) = a\lambda(f) + b\lambda(g), \quad (23.13)$$

whenever $a, b \geq 0$ and $f, g \geq 0$. Then there exists a unique extension of λ to all of $C_0(X)$, Λ such that whenever $f, g \in C_0(X)$ and $a, b \in \mathbb{C}$, it follows $\Lambda(af + bg) = a\Lambda(f) + b\Lambda(g)$. If

$$|\lambda(f)| \leq C\|f\|_\infty \quad (23.14)$$

then $|\Lambda f| \leq \lambda(|f|) \leq C\|f\|_\infty$.

Proof: There is only one possible way to extend this functional to obtain a linear functional and the arguments are identical with those of Theorem 10.7.8 so I will refer to this earlier theorem for these arguments. In particular, you must have

$$\begin{aligned}\Lambda(f) &= \Lambda(\operatorname{Re} f) + i\Lambda(\operatorname{Im} f) = \Lambda(\operatorname{Re} f)^+ - \Lambda(\operatorname{Re} f)^- \\ &\quad + i(\Lambda(\operatorname{Im} f)^+ - \Lambda(\operatorname{Im} f)^-) \\ &= \lambda(\operatorname{Re} f)^+ - \lambda(\operatorname{Re} f)^- + i(\lambda(\operatorname{Im} f)^+ - \lambda(\operatorname{Im} f)^-)\end{aligned}$$

Since the nature of the functions is different, being continuous here rather than only measurable, the only thing left is to show the claim about continuity of Λ in case of 23.14.

What of the last claim that $|\Lambda f| \leq \lambda(|f|)$? Let ω have $|\omega| = 1$ and $|\Lambda f| = \omega \Lambda(f)$. Since Λ is linear,

$$\begin{aligned}|\Lambda f| &= \omega \Lambda(f) = \Lambda(\omega f) = \Lambda(\operatorname{Re} \omega f) \leq \Lambda(\operatorname{Re}(\omega f)^+) \\ &= \lambda(\operatorname{Re}(\omega f)^+) \leq \lambda(|f|) \leq C \|f\|_\infty \blacksquare\end{aligned}$$

Corollary 23.7.5 *Let $L \in \mathcal{L}(C_0(X), \mathbb{C})$. Then there exists $\Lambda \in \mathcal{L}(C_0(X), \mathbb{C})$ which satisfies $\|\Lambda\| = \|L\|$ but also Λ is a positive linear functional meaning if $f \geq 0$, then $\Lambda(f) \geq 0$.*

Proof: Let λ be the righteous functional defined in Lemma 23.7.3 which satisfies $|\lambda(f)| \leq \|L\| \|f\|_\infty$. Then let Λ be its extension defined in Lemma 23.7.4 which also satisfies $|\Lambda(f)| \leq \|L\| \|f\|_\infty$. Then this is a positive linear functional and $\|\Lambda\| \leq \|L\|$. However, from the definition of λ ,

$$|Lg| \leq \lambda(|g|) = \Lambda(|g|) \leq \|\Lambda\| \|g\|_\infty$$

and so also $\|L\| \leq \|\Lambda\|$. ■

23.7.2 The Riesz Representation Theorem

What follows is the Riesz representation theorem for $C_0(X)'$.

Theorem 23.7.6 *Let $L \in (C_0(X))'$ for X a locally compact Hausdorff space. Then there exists a σ algebra \mathcal{F} and a finite Radon measure μ and a function $\sigma \in L^\infty(X, \mu)$ such that for all $f \in C_0(X)$,*

$$L(f) = \int_X f \sigma d\mu.$$

Furthermore, $\mu(X) = \|L\|$, $|\sigma| = 1$ a.e. and if $v(E) \equiv \int_E \sigma d\mu$ then $\mu = |v|$.

Proof: From Corollary 23.7.5, there exists a positive linear functional Λ defined on $C_0(X)$ with $\|\Lambda\| = \|L\|$. Then let μ be the Radon measure representing Λ for which, by Lemma 23.7.2, $\mu(X) = \|\Lambda\| = \|L\|$.

For $f \in C_c(X)$, $|Lf| \leq \lambda(|f|) = \Lambda(|f|) = \int_X |f| d\mu = \|f\|_{L^1(\mu)}$. Since μ is both inner and outer regular thanks to it being finite, $C_c(X)$ is dense in $L^1(X, \mu)$. (See Theorem 12.2.4 for more than is needed.) Thus there is a unique extension of L to $\tilde{L} \in (L^1(X, \mu))'$ and

by the Riesz representation theorem for the dual of $L^1(\mu)$, there exists $\sigma \in L^\infty(\mu)$ with $\tilde{L}(f) = \int_X f \sigma d\mu$. In particular,

$$L(f) = \int_X f \sigma d\mu \text{ for } f \in C_0(X).$$

It remains to verify that $|\sigma| = 1$.

If E is measurable, the regularity of μ implies there exists a sequence of nonnegative bounded functions $f_n \in C_c(X)$ such that $f_n(x) \rightarrow \chi_E(x)$ a.e. and in $L^1(\mu)$. Then using the dominated convergence theorem,

$$\int_E d\mu = \lim_{n \rightarrow \infty} \int_X f_n d\mu = \lim_{n \rightarrow \infty} \Lambda(f_n) \geq \lim_{n \rightarrow \infty} |L f_n| = \lim_{n \rightarrow \infty} \left| \int_X f_n \sigma d\mu \right| = \left| \int_E \sigma d\mu \right|$$

and so if $\mu(E) > 0$, $1 \geq \left| \frac{1}{\mu(E)} \int_E \sigma d\mu \right|$ which shows from Lemma 23.2.7 that $|\sigma| \leq 1$ a.e.

But also, choosing f_1 appropriately, $\|f_1\|_\infty \leq 1$, $|L f_1| + \varepsilon > \|L\| = \mu(X)$. Letting $\omega(L f_1) = |L f_1|$, $|\omega| = 1$,

$$\begin{aligned} \mu(X) &= \|L\| = \sup_{\|f\|_\infty \leq 1} |L f| \leq |L f_1| + \varepsilon = \omega L f_1 + \varepsilon = \int_X f_1 \omega \sigma d\mu + \varepsilon \\ &= \int_X \operatorname{Re}(f_1 \omega \sigma) d\mu + \varepsilon \leq \int_X |\sigma| d\mu + \varepsilon \leq \mu(X) + \varepsilon \end{aligned}$$

and since ε is arbitrary, $\mu(X) \leq \int_X |\sigma| d\mu \leq \mu(X)$ which requires $|\sigma| = 1$ a.e. since it was shown to be no larger than 1 and if it is smaller than 1 on a set of positive measure, then the above could not hold.

If $\nu(E) \equiv \int_E \sigma d\mu$, by Corollary 23.2.9, $|\nu|(E) = \int_E |\sigma| d\mu = \int_E 1 d\mu = \mu(E)$ ■

Sometimes people write $L(f) = \int_X f d\nu \equiv \int_X f \sigma d\mu$ where $\sigma d\mu$ is the polar decomposition of the complex measure $\nu(E) \equiv \int_E \sigma d\mu$. Then with this convention, the above representation is $L(f) = \int_X f d\nu$, $|\nu|(X) = \|L\|$. Also note that at most one ν can represent L . If there were two of them $\nu_i, i = 1, 2$, then $0 = \int_X f \sigma d|\nu_1 - \nu_2|$, $|\sigma| = 1$, and so it will follow that $|\nu_1 - \nu_2|(X) = 0$ because you could approximate $\bar{\sigma}$ with a sequence f_n and after using the dominated convergence theorem, you would get $|\nu_1 - \nu_2|(X) = 0$. Hence $\nu_1 = \nu_2$.

Corollary 23.7.7 *Let $L \in \mathcal{L}(C_0(X), \mathbb{C})$. Then there exists a unique complex measure ν such that for all $f \in C_0(X)$, $L(f) = \int_X f d\nu$ and $|\nu|(X) = \|L\|$.*

23.8 Exercises

1. Suppose μ is a vector measure having values in \mathbb{R}^n or \mathbb{C}^n . Can you show that $|\mu|$ must be finite? **Hint:** You might define for each e_i , one of the standard basis vectors, the real or complex measure, μ_{e_i} given by $\mu_{e_i}(E) \equiv e_i \cdot \mu(E)$. Why would this approach not yield anything for an infinite dimensional normed linear space in place of \mathbb{R}^n ? Have a look at the proof of Theorem 23.1.3.
2. The Riesz representation theorem of the L^p spaces can be used to prove a very interesting inequality. Let $r, p, q \in (1, \infty)$ satisfy $\frac{1}{r} = \frac{1}{p} + \frac{1}{q} - 1$. Then $\frac{1}{q} = 1 + \frac{1}{r} - \frac{1}{p} > \frac{1}{r}$

and so $r > q$. Let $\theta \in (0, 1)$ be chosen so that $\theta r = q$. Then also $\frac{1}{r} = \left(\overbrace{1 - \frac{1}{p'}}^{1/p + 1/p' = 1} \right) + \frac{1}{q} - 1 = \frac{1}{q} - \frac{1}{p'}$ and so $\frac{\theta}{q} = \frac{1}{q} - \frac{1}{p'}$ which implies $p'(1 - \theta) = q$. Now let $f \in L^p(\mathbb{R}^n)$, $g \in L^q(\mathbb{R}^n)$, $f, g \geq 0$. Justify the steps in the following argument using what was just shown that $\theta r = q$ and $p'(1 - \theta) = q$. Let $h \in L^{r'}(\mathbb{R}^n)$. ($\frac{1}{r} + \frac{1}{r'} = 1$),

$$\begin{aligned}
& \left| \int f * g(x) h(x) dx \right| \\
&= \left| \int \int f(y) g(x - y) h(x) dx dy \right| \\
&\leq \int \int |f(y)| |g(x - y)|^\theta |g(x - y)|^{1-\theta} |h(x)| dy dx \\
&\leq \int \left(\int (|g(x - y)|^{1-\theta} |h(x)|)^{r'} dx \right)^{1/r'} \\
&\quad \left(\int (|f(y)| |g(x - y)|^\theta)^r dx \right)^{1/r} dy \\
&\leq \left[\int \left(\int (|g(x - y)|^{1-\theta} |h(x)|)^{r'} dx \right)^{p'/r'} dy \right]^{1/p'} \\
&\quad \left[\int \left(\int (|f(y)| |g(x - y)|^\theta)^r dx \right)^{p/r} dy \right]^{1/p} \\
&\leq \left[\int \left(\int (|g(x - y)|^{1-\theta} |h(x)|)^{p'} dy \right)^{r'/p'} dx \right]^{1/r'} \\
&\quad \left[\int |f(y)|^p \left(\int |g(x - y)|^{\theta r} dx \right)^{p/r} dy \right]^{1/p} \\
&= \left[\int |h(x)|^{r'} \left(\int |g(x - y)|^{(1-\theta)p'} dy \right)^{r'/p'} dx \right]^{1/r'} \|g\|_q^{q/r} \|f\|_p \\
&= \|g\|_q^{q/r} \|g\|_q^{q/p'} \|f\|_p \|h\|_{r'} = \|g\|_q \|f\|_p \|h\|_{r'}. \tag{23.15}
\end{aligned}$$

Young's inequality says that

$$\|f * g\|_r \leq \|g\|_q \|f\|_p. \tag{23.16}$$

Therefore $\|f * g\|_r \leq \|g\|_q \|f\|_p$. How does this inequality follow from the above computation? Does 23.15 continue to hold if r, p, q are only assumed to be in $[1, \infty]$? Explain. Does 23.16 hold even if r, p , and q are only assumed to lie in $[1, \infty]$?

3. Suppose $(\Omega, \mu, \mathcal{F})$ is a finite measure space and that $\{f_n\}$ is a sequence of functions which converge weakly to 0 in $L^p(\Omega)$. This means that $\int_{\Omega} f_n g d\mu \rightarrow 0$ for every $g \in L^{p'}(\Omega)$. Suppose also that $f_n(x) \rightarrow 0$ a.e. Show that then $f_n \rightarrow 0$ in $L^{p-\varepsilon}(\Omega)$ for every $p > \varepsilon > 0$.
4. Give an example of a sequence of functions in $L^\infty(-\pi, \pi)$ which converges weak * to zero but which does not converge pointwise a.e. to zero. Convergence weak * to 0 means that for every $g \in L^1(-\pi, \pi)$, $\int_{-\pi}^{\pi} g(t) f_n(t) dt \rightarrow 0$. **Hint:** First consider $g \in C_c^\infty(-\pi, \pi)$ and maybe try something like $f_n(t) = \sin(nt)$. Do integration by parts.
5. Let (Ω, \mathcal{F}) be a measurable space and let $\lambda : \mathcal{F} \rightarrow (-\infty, \infty]$ be such that if the E_i are disjoint sets in \mathcal{F} then $\lambda(\cup_i E_i) = \sum_i \lambda(E_i)$ where this sum either equals a real number or $+\infty$. The Hahn decomposition says there exist measurable sets P, N such that $P \cup N = \Omega$, $P \cap N = \emptyset$, and for each $F \subseteq P$, $\lambda(F) \geq 0$ and for each $F \subseteq N$, $\lambda(F) \leq 0$. These sets P, N are called the positive set and the negative set respectively. Show the existence of the Hahn decomposition. Also explain how this decomposition is unique in the sense that if P', N' is another Hahn decomposition, then $(P \setminus P') \cup (P' \setminus P)$ has measure zero, a similar formula holding for N, N' . When you have the Hahn decomposition, as just described, you define $\lambda^+(E) \equiv \lambda(E \cap P)$, $\lambda^-(E) \equiv -\lambda(E \cap N)$. This is sometimes called the Hahn Jordan decomposition. **Hint:** You could use similar arguments leading to Theorem 10.13.5. However, this time be sure that the Hausdorff maximality argument is applied to sets which have negative measure.
6. From Problem 5 for λ having values in $(-\infty, \infty]$ you have the Hahn Jordan decomposition for the measure λ , $\lambda^+(E) \equiv \lambda(E \cap P)$, $\lambda^-(E) \equiv -\lambda(E \cap N)$. Explain why λ^- is a finite measure. **Hint:** It is a complex measure which happens to have values in \mathbb{R} .
7. If $\mu : \mathcal{F} \rightarrow [0, \infty)$ is a finite measure and if $\lambda : \mathcal{F} \rightarrow (-\infty, \infty)$ is another signed measure and $\lambda \ll \mu$ meaning that if $\mu(E) = 0$, then $\lambda(E) = 0$, show that $\lambda^+, \lambda^- \ll \mu$ with both being finite measures. Explain why there exists $f \in L^1(\Omega)$ with $\lambda(E) = \int \mathcal{X}_E f d\mu$.
8. Suppose λ is like the above problem but has values in $[0, \infty]$ and μ is a finite real valued measure on \mathcal{F} . Suppose also that $\lambda \ll \mu$. Show there exists a measurable $f \geq 0$ such that $\lambda(E) = \int \mathcal{X}_E f d\mu$.
9. What if λ has values in $[-\infty, \infty)$. Prove there exists a Hahn decomposition for λ as in the above problem. Why do we not allow λ to have values in $[-\infty, \infty]$? **Hint:** You might want to consider $-\lambda$.
10. Suppose X is a Banach space and let X' denote its dual space. A sequence $\{x_n^*\}_{n=1}^\infty$ in X' is said to converge weak * to $x^* \in X'$ if for every $x \in X$, $\lim_{n \rightarrow \infty} x_n^*(x) = x^*(x)$. Let $\{\phi_n\}$ be a mollifier. Also let δ be the measure defined by $\delta(E) = 1$ if $0 \in E$ and 0 if $1 \notin E$. Explain how $\phi_n \rightarrow \delta$ weak *.
11. It was shown above that if $\phi \in X'$ where X is a uniformly convex Banach space, then there exists $x \in X$, $\|x\| = 1$, and $\phi(x) = \|\phi\|$. Show that this x must be unique. **Hint:** Recall that uniform convexity implies strict convexity.

12. Suppose $\lambda(E) = \int_E h d\mu$ where h is real valued and μ is a finite measure so that λ is also real valued. Let P, N be a Hahn decomposition for λ . Show that $|\lambda|(E) = \int_E |h| d\mu$. **Hint:** Argue that on P it follows $h \geq 0$ a.e. and on $N, h \leq 0$ a.e. Then estimate $\sum_{F \in \pi(E)} \lambda(F)$ using a Hahn decomposition. If we defined $|x + iy|_1$ as $|x|_1 + |y|_1$, and the total variation exactly the same way for a complex valued measure except for letting $|\cdot|_1$ refer to this way of measuring magnitude, then everything would be much easier. Why don't we do this and save a lot of trouble?

Chapter 24

The Bochner Integral

From my experience, the Bochner integral tends to be ignored. However, it is one of the most useful and important ideas in functional analysis, at least in my experience. Perhaps it is not useful in number theory or modern algebra, but I have used it in almost every paper I have written during my career. Therefore, I am including it in this book. If people are not interested in it, they can ignore it, but if so, they will be missing out on some very nice mathematics. The work of Pettis about weak and strong measurability is particularly interesting.

24.1 Strong and Weak Measurability

In this section (Ω, \mathcal{F}) will be a measurable space and X will be a Banach space which contains the values of either a function or a measure. The Banach space will be either a real or a complex Banach space but the field of scalars does not matter and so it is denoted by \mathbb{F} with the understanding that $\mathbb{F} = \mathbb{C}$ unless otherwise stated. The theory presented here includes the case where $X = \mathbb{R}^n$ or \mathbb{C}^n but it does not include the situation where f could have values in something like $[0, \infty]$ which is not a vector space. To begin with here is a definition.

Definition 24.1.1 A function, $x : \Omega \rightarrow X$, for X a Banach space, is finitely valued and measurable if it is of the form $x(\omega) = \sum_{i=1}^n a_i \chi_{B_i}(\omega)$ where $B_i \in \mathcal{F}$ for each i . These are called simple functions. A function x from Ω to X is said to be strongly measurable if there exists a sequence of finitely valued and measurable functions $\{x_n\}$ with $x_n(\omega) \rightarrow x(\omega)$. The function x is said to be weakly measurable if, for each $f \in X'$, $f \circ x$ is a scalar valued measurable function.

The approximating simple functions can be modified so that the norm of each is no more than $2 \|x(\omega)\|$. This is a useful observation.

Lemma 24.1.2 Let x be strongly measurable. Then $\|x\|$ is a real valued measurable function. There exists a sequence of simple functions $\{y_n\}$ which converges to $x(\omega)$ pointwise and also $\|y_n(\omega)\| \leq 2 \|x(\omega)\|$ for all ω .

Proof: Consider the first claim. Letting x_n be a sequence of simple functions converging to x pointwise, it follows that $\|x_n\|$ is a real valued measurable function. Since $\|x\|$ is a pointwise limit, $\|x\|$ is a real valued measurable function.

Let $\lim_{n \rightarrow \infty} x_n(\omega) = x(\omega)$ where $x_n(\omega) \equiv \sum_{k=1}^{m_n} a_k^n \chi_{E_k^n}(\omega)$. Then define

$$y_n(\omega) \equiv \begin{cases} x_n(\omega) & \text{if } \|x_n(\omega)\| < 2 \|x(\omega)\| \\ 0 & \text{if } \|x_n(\omega)\| \geq 2 \|x(\omega)\| \end{cases}$$

so $y_n(\omega) = \sum_{k=1}^{m_n} a_k^n \chi_{E_k^n \cap [\|a_k^n\| \leq 2\|x\|]}(\omega)$. It follows y_n is a simple function. If $\|x(\omega)\| = 0$, then $y_n(\omega) = 0$ and so $y_n(\omega) \rightarrow x(\omega)$. If $\|x(\omega)\| > 0$, then eventually, $y_n(\omega) = x_n(\omega)$ and so in this case, $y_n(\omega) \rightarrow x(\omega)$. ■

Earlier, a function was measurable if inverse images of open sets were measurable. Something similar holds here. The difference is that another condition needs to hold about the values being separable. First is a somewhat obvious lemma.

Lemma 24.1.3 Suppose S is a nonempty subset of a metric space (X, d) and $S \subseteq T$ where T is separable. Then there exists a countable dense subset of S .

Proof: Let D be the countable dense subset of T . Now consider the countable set \mathcal{B} of balls having center at a point of D and radius a positive rational number such that also, each ball in \mathcal{B} has nonempty intersection with S . Let \mathcal{D} consist of a point from $S \cap B$ whenever $B \in \mathcal{B}$ (axiom of choice). Let $s \in S$ and consider $B(s, \varepsilon)$. Let r be rational with $r < \varepsilon$. Now $B(s, \frac{r}{10})$ contains a point $d \in D$. Thus $B(d, \frac{r}{10}) \in \mathcal{B}$ and $s \in B(d, \frac{r}{10})$. Let $\hat{d} \in \mathcal{D} \cap B(d, \frac{r}{10})$. Thus $d(s, \hat{d}) < \frac{r}{5} < r < \varepsilon$ so $\hat{d} \in B(s, \varepsilon)$ and this shows that \mathcal{D} is a countable dense subset of S as claimed. ■

The following is a general result in metric space.

Lemma 24.1.4 Let X be a metric space and suppose V is a nonempty open set. Then there exists open sets V_m such that

$$\cdots V_m \subseteq \overline{V_m} \subseteq V_{m+1} \subseteq \cdots, \quad V = \bigcup_{m=1}^{\infty} V_m. \quad (24.1)$$

Proof: Recall that if S is a nonempty set, $x \rightarrow \text{dist}(x, S)$ is a continuous map from X to \mathbb{R} . First assume $V \neq X$. Let $V_m \equiv \{x \in V : \text{dist}(x, V^C) > \frac{1}{m}\}$. Then for large enough m , this set is nonempty and contained in V . Furthermore, if $x \in V$ then it is at a positive distance to the closed set V^C so eventually, $x \in V_m$. Now

$$V_m \subseteq \overline{V_m} \subseteq \left\{x \in V : \text{dist}(x, V^C) \geq \frac{1}{m}\right\} \subseteq V_{m+1} \subseteq V$$

Indeed, if p is a limit point of V_m , then there are $x_n \in V_m$ with $x_n \rightarrow p$. Thus $\text{dist}(x_n, V^C) \rightarrow \text{dist}(p, V^C)$ and so p is in $\{x \in V : \text{dist}(x, V^C) \geq \frac{1}{m}\}$. ■

Theorem 24.1.5 x is strongly measurable if and only if $x^{-1}(U)$ is measurable for all U open in X and $x(\Omega)$ is separable. Thus, if X is separable, x is strongly measurable if and only if $x^{-1}(U)$ is measurable for all U open.

Proof: \Leftarrow Suppose first $x^{-1}(U)$ is measurable for all U open in X and $x(\Omega)$ is separable. It follows $x^{-1}(B)$ is measurable for all B Borel because $\{B : x^{-1}(B) \text{ is measurable}\}$ is a σ algebra containing the open sets. Let $\{a_n\}_{n=1}^{\infty}$ be the dense subset of $x(\Omega)$. Let

$$U_k^n \equiv \{z \in X : \|z - a_k\| \leq \min\{\|z - a_l\|\}_{l=1}^n\}.$$

In words, U_k^n is the set of points of X which are as close to a_k as they are to any of the a_l for $l \leq n$.

$$B_k^n \equiv x^{-1}(U_k^n), \quad D_k^n \equiv B_k^n \setminus \left(\bigcup_{i=1}^{k-1} B_i^n\right), \quad D_1^n \equiv B_1^n,$$

and $x_n(\omega) \equiv \sum_{k=1}^n a_k \mathcal{X}_{D_k^n}(\omega)$. Thus $x_n(\omega)$ is a closest approximation to $x(\omega)$ from $\{a_k\}_{k=1}^n$ and so $x_n(\omega) \rightarrow x(\omega)$ because $\{a_n\}_{n=1}^{\infty}$ is dense in $x(\Omega)$. Furthermore, x_n is measurable because each D_k^n is measurable.

\Rightarrow Now suppose that x is strongly measurable. Then some sequence of measurable finite valued functions $\{x_n\}$ converges pointwise to x . Then $x_n^{-1}(W)$ is measurable for every open set W because it is just a finite union of measurable sets. If $x_n(\omega) = \sum_{k=1}^n c_k \mathcal{X}_{E_k}(\omega)$,

then $x_n^{-1}(W) = \cup \{E_k : c_k \in W\}$. Thus, $x_n^{-1}(W)$ is measurable for every Borel set W . This follows from the observation that $\{W : x_n^{-1}(W) \text{ is measurable}\}$ is a σ algebra containing the open sets. Since X is a metric space, it follows that if U is an open set in X , there exists a sequence of open sets, $\{V_n\}$ which satisfies

$$\bar{V}_n \subseteq U, \bar{V}_n \subseteq V_{n+1}, U = \cup_{n=1}^{\infty} V_n.$$

Then $x^{-1}(V_m) \subseteq \bigcup_{n < \infty} \bigcap_{k \geq n} x_k^{-1}(V_m) \subseteq x^{-1}(\bar{V}_m)$. This implies

$$x^{-1}(U) = \bigcup_{m < \infty} x^{-1}(V_m) \subseteq \bigcup_{m < \infty} \bigcup_{n < \infty} \bigcap_{k \geq n} x_k^{-1}(V_m) \subseteq \bigcup_{m < \infty} x^{-1}(\bar{V}_m) \subseteq x^{-1}(U).$$

Since $x^{-1}(U) = \bigcup_{m < \infty} \bigcup_{n < \infty} \bigcap_{k \geq n} x_k^{-1}(V_m)$, it follows that $x^{-1}(U)$ is measurable for every open U . It remains to show $x(\Omega)$ is separable. Let $D \equiv$ all values of the x_n . Then $x(\Omega) \subseteq \bar{D}$, which has a countable dense subset. By Lemma 24.1.3, $x(\Omega)$ is separable. ■

Lemma 24.1.6 *Let $x \in X$ a normed linear space. Then there exists $f \in X'$ such that $\|f\| = 1$ and $f(x) = \|x\|$.*

Proof: Consider the one dimensional subspace $M \equiv \left\{ \alpha \frac{x}{\|x\|} : \alpha \in \mathbb{R} \right\}$ and define a continuous linear functional on M by $g\left(\alpha \frac{x}{\|x\|}\right) \equiv \alpha$. Then the operator norm of g is obtained as $\|g\| \equiv \sup_{|\alpha| \leq 1} |\alpha| = 1$. Extend g to all of X using the Hahn Banach theorem, calling the extended function f . Then $\|f\| = 1$ and $f(x) = f\left(\|x\| \frac{x}{\|x\|}\right) \equiv \|x\|$. ■

The next lemma is interesting for its own sake. Roughly it says that if a Banach space is separable, then the unit ball in the dual space is weak * separable. This will be used to prove Pettis's theorem, one of the major theorems in this subject which relates weak measurability to strong measurability. First here is a standard application which comes from earlier material on the Hahn Banach theorem.

Lemma 24.1.7 *If X is a separable Banach space with B' the closed unit ball in X' , then there exists a sequence $\{f_n\}_{n=1}^{\infty} \equiv D' \subseteq B'$ with the property that for every $x \in X$, $\|x\|$ is obtained as $\|x\| = \sup_{f \in D'} |f(x)|$. If H is a dense subset of X' then D' may be chosen to be contained in H .*

Proof: Let $\{a_k\}_{k=1}^{\infty}$ be a countable dense set in X . Consider the mapping $\phi_n : B' \rightarrow \mathbb{R}^n$ given by $\phi_n(f) \equiv (f(a_1), \dots, f(a_n))$.

Then $\phi_n(B')$ is contained in a compact subset of \mathbb{R}^n because $|f(a_k)| \leq \|a_k\|$. Therefore, there exists a countable dense subset of $\phi_n(B')$, $\{\phi_n(f_k)\}_{k=1}^{\infty}$. Pick $h_j^k \in H \cap B'$ such that $\lim_{j \rightarrow \infty} \|f_k - h_j^k\| = 0$. Then $\left\{ \phi_n(h_j^k) \right\}_{k,j}$ must also be dense in $\phi_n(B')$. Let $D'_n = \left\{ h_j^k \right\}_{k,j}$. Thus D'_n is a countable collection of $f \in B'$ which can be used to approximate each $\|a_k\|$, $k \leq n$. Indeed, if x is arbitrary, there exists $f_x \in B'$ with $f_x(x) = \|x\|$ and so if $x = a_k$, then $\|a_k\|$ will be close to $g(a_k)$ for some $g \in D'_n$. Define $D' \equiv \cup_{n=1}^{\infty} D'_n$.

From the construction, D' is countable and can be used to approximate each $\|a_m\|$. That is, $\|a_m\| = \sup \{|f(a_m)| : f \in D'\}$. Then, for x arbitrary, $|f(x)| \leq \|x\|$ and so

$$\begin{aligned} \|x\| &\leq \|x - a_m\| + \|a_m\| = \|x - a_m\| + \sup \{|f(a_m)| : f \in D'\} \\ &\leq \|x - a_m\| + \sup \{|f(a_m - x) + f(x)| : f \in D'\} \\ &\leq \sup \{|f(x)| : f \in D'\} + 2\|x - a_m\| \leq \|x\| + 2\|x - a_m\|. \end{aligned}$$

Since a_m is arbitrary and the $\{a_m\}_{m=1}^\infty$ are dense, this establishes the claim of the lemma. ■

Note that the proof would work the same if H were only given to be weak $*$ dense.

The next theorem is one of the most important results in the subject. It is due to Pettis and appeared in 1938 [45].

Theorem 24.1.8 *If x has values in a separable Banach space X , then x is weakly measurable if and only if x is strongly measurable.*

Proof: \Rightarrow It is necessary to show $x^{-1}(U)$ is measurable whenever U is open. Since every open set is a countable union of balls, it suffices to show $x^{-1}(B(a, r))$ is measurable for any ball $B(a, r)$. Since, $B(x, r) = \bigcup_{n=1}^\infty \overline{B(x, (1 - \frac{1}{n})r)}$ or by Lemma 24.1.4, every open ball is the countable union of closed balls, it suffices to verify $x^{-1}(\overline{B(a, r)})$ is measurable. For D' described in Lemma 24.1.7,

$$\begin{aligned} x^{-1}(\overline{B(a, r)}) &= \{\omega : \|x(\omega) - a\| \leq r\} = \left\{ \omega : \sup_{f \in D'} |f(x(\omega) - a)| \leq r \right\} \\ &= \bigcap_{f \in D'} \{\omega : |f(x(\omega) - a)| \leq r\} \\ &= \bigcap_{f \in D'} \{\omega : |f(x(\omega)) - f(a)| \leq r\} \\ &= \bigcap_{f \in D'} (f \circ x)^{-1} \overline{B(f(a), r)} \end{aligned}$$

which equals a countable intersection of measurable sets because it is assumed that $f \circ x$ is measurable for all $f \in X'$.

\Leftarrow Next suppose x is strongly measurable. Then there exists a sequence of simple functions x_n which converges to x pointwise. Hence for all $f \in X'$, $f \circ x_n$ is measurable since f is continuous and $f \circ x_n \rightarrow f \circ x$ pointwise. Thus x is weakly measurable. ■

The same method of proof yields the following interesting corollary.

Corollary 24.1.9 *Let X be a separable Banach space and let $\mathcal{B}(X)$ denote the σ algebra of Borel sets. Let H be a dense subset of X' . Then $\mathcal{B}(X) = \sigma(H) \equiv \mathcal{F}$, where $\sigma(H)$ is the smallest σ algebra of subsets of X which has the property that every function, $x^* \in H$ is measurable. That is $(x^*)^{-1}(\text{open}) \in \mathcal{F}$.*

Proof: First I need to show \mathcal{F} contains open balls because then \mathcal{F} will contain the open sets and hence the Borel sets. As noted above, it suffices to show \mathcal{F} contains closed balls. Let D' be those functionals in B' defined in Lemma 24.1.7 contained in H . Then

$$\begin{aligned} \{x : \|x - a\| \leq r\} &= \left\{ x : \sup_{x^* \in D'} |x^*(x - a)| \leq r \right\} \\ &= \bigcap_{x^* \in D'} \{x : |x^*(x - a)| \leq r\} \\ &= \bigcap_{x^* \in D'} \{x : |x^*(x) - x^*(a)| \leq r\} \\ &= \bigcap_{x^* \in D'} x^{*-1}(\overline{B(x^*(a), r)}) \in \sigma(H) \end{aligned}$$

which is measurable because this is a countable intersection of measurable sets. Thus \mathcal{F} contains closed balls, hence open balls, hence open sets so $\sigma(H) \equiv \mathcal{F} \supseteq \mathcal{B}(X)$.

To show the other direction for the inclusion, note that each x^* is $\mathcal{B}(X)$ measurable because $x^{*-1}(\text{open set}) = \text{open set}$. Therefore, $\mathcal{B}(X) \supseteq \sigma(H)$. ■

What of limits of measurable functions? The next theorem says that the usual theorem about limits of measurable functions being measurable holds. The proof is similar to showing that the limit of measurable finitely valued functions is measurable given above.

Theorem 24.1.10 *Let x_n and x be functions mapping Ω to X where \mathcal{F} is a σ -algebra of measurable sets of Ω and X is a Banach space. Thus X satisfies 24.1. Then if x_n is strongly measurable, and $x(\omega) = \lim_{n \rightarrow \infty} x_n(\omega)$, it follows that x is also strongly measurable. (Pointwise limits of measurable functions are measurable.)*

Proof: Let $\{V_m\}$ be the sequence of 24.1. Since x is the pointwise limit of x_n ,

$$x^{-1}(V_m) \subseteq \{\omega : x_k(\omega) \in V_m \text{ for all } k \text{ large enough}\} \subseteq x^{-1}(\overline{V_m}).$$

Therefore,

$$\begin{aligned} x^{-1}(V) &= \bigcup_{m=1}^{\infty} x^{-1}(V_m) \subseteq \bigcup_{m=1}^{\infty} \bigcup_{n=1}^{\infty} \bigcap_{k=n}^{\infty} x_k^{-1}(V_m) \\ &\subseteq \bigcup_{m=1}^{\infty} x^{-1}(\overline{V_m}) = x^{-1}(V). \end{aligned}$$

It follows $x^{-1}(V) \in \mathcal{F}$ because it equals the expression in the middle which is measurable. Note that this shows the characterization of measurability in terms of inverse images of open sets being measurable sets. Thus the theorem is proved in the case of separable Banach spaces. However, Lemma 24.1.3 can be applied to conclude that this holds in general because each x_n is separably valued given they are each strongly measurable and $x(\Omega) \subseteq \overline{D}$ where $D = \bigcup_n D_n$ for D_n a countable dense subset of $x_n(\Omega)$. ■

Note that the same conclusion in terms of inverse images being measurable would hold for any metric space.

Corollary 24.1.11 *x is strongly measurable if and only if $x(\Omega)$ is separable and x is weakly measurable.*

Proof: Strong measurability clearly implies weak measurability. If $x_n(\omega) \rightarrow x(\omega)$ where x_n is simple, then $f(x_n(\omega)) \rightarrow f(x(\omega))$ for all $f \in X'$. Hence $f \circ x$ is measurable by Theorem 24.1.10 because it is the limit of a sequence of measurable functions. Let D denote the set of all values of the x_n . Then \overline{D} is a separable set containing $x(\Omega)$. Thus \overline{D} is a separable metric space. Therefore $x(\Omega)$ is separable also by the last part of the proof of Theorem 24.1.5.

Now suppose D is a countable dense subset of $x(\Omega)$ and x is weakly measurable. Let Z be the subset consisting of all finite linear combinations of D with the scalars coming from the set of rational points of \mathbb{F} . Thus, Z is countable. Letting $Y = \overline{Z}$, Y is a separable Banach space containing $x(\Omega)$. If $f \in Y'$, f can be extended to an element of X' by the Hahn Banach theorem. Therefore, x is a weakly measurable Y valued function. Now use Theorem 24.1.8 to conclude x is strongly measurable. ■

Weakly measurable as defined above means $\omega \rightarrow x^*(x(\omega))$ is measurable for every $x^* \in X'$. The next lemma ties this weak measurability to the usual version of measurability in which a function is measurable when inverse images of open sets are measurable.

Lemma 24.1.12 *Let X be a Banach space and let $x : (\Omega, \mathcal{F}) \rightarrow K \subseteq X$ where K is weakly compact and X' is separable. Then x is weakly measurable if and only if $x^{-1}(U) \in \mathcal{F}$ whenever U is a weakly open set.*

Proof: By Corollary 21.5.11 on Page 559, there exists a metric d , such that the metric space topology with respect to d coincides with the weak topology on K . Since K is compact, it follows that K is also separable. Hence it is completely separable and so there exists a countable basis of open sets \mathcal{B} for the weak topology on K . It follows that if U is any weakly open set, covered by basic sets of the form $B_A(x, r)$ where A is a finite subset of X' , there exists a countable collection of these sets of the form $B_A(x, r)$ which covers U .

Suppose now that x is weakly measurable. To show $x^{-1}(U) \in \mathcal{F}$ whenever U is weakly open, it suffices to verify $x^{-1}(B_A(z, r)) \in \mathcal{F}$ for any set, $B_A(z, r)$. Let $A = \{x_1^*, \dots, x_m^*\}$. Then

$$\begin{aligned} x^{-1}(B_A(z, r)) &= \{\omega \in \Omega : \rho_A(x(\omega) - z) < r\} \\ &\equiv \left\{ \omega \in \Omega : \max_{x^* \in A} |x^*(x(\omega) - z)| < r \right\} \\ &= \bigcup_{i=1}^m \{\omega \in \Omega : |x_i^*(x(\omega) - z)| < r\} \\ &= \bigcup_{i=1}^m \{\omega \in \Omega : |x_i^*(x(\omega)) - x_i^*(z)| < r\} \end{aligned}$$

which is measurable because each $x_i^* \circ x$ is given to be measurable.

Next suppose $x^{-1}(U) \in \mathcal{F}$ whenever U is weakly open. Then in particular this holds when $U = B_{x^*}(z, r)$ for arbitrary x^* . Hence

$$\{\omega \in \Omega : x(\omega) \in B_{x^*}(z, r)\} \in \mathcal{F}.$$

But this says the same as

$$\{\omega \in \Omega : |x^*(x(\omega)) - x^*(z)| < r\} \in \mathcal{F}$$

Since $x^*(z)$ can be a completely arbitrary element of \mathbb{F} , it follows $x^* \circ x$ is an \mathbb{F} valued measurable function. In other words, x is weakly measurable according to the former definition. ■

One can also define weak $*$ measurability and prove a theorem just like the Pettis theorem above. The next lemma is the analogue of Lemma 24.1.7.

Lemma 24.1.13 *Let B be the closed unit ball in X . If X' is separable, there exists a sequence $\{x_m\}_{m=1}^\infty \equiv D \subseteq B$ with the property that for all $y^* \in X'$, $\|y^*\| = \sup_{x \in D} |y^*(x)|$.*

Proof: Let $\{x_k^*\}_{k=1}^\infty$ be the dense subset of X' . Define $\phi_n : B \rightarrow \mathbb{F}^n$ by the formula $\phi_n(x) \equiv (x_1^*(x), \dots, x_n^*(x))$.

Then $|x_k^*(x)| \leq \|x_k^*\|$ and so $\phi_n(B)$ is contained in a compact subset of \mathbb{F}^n . Therefore, there exists a countable set, $D_n \subseteq B$ such that $\phi_n(D_n)$ is dense in $\phi_n(B)$. That is, $\{(x_1^*(x), \dots, x_n^*(x)) : x \in D_n\}$ is dense in $\phi_n(B)$. $D \equiv \bigcup_{n=1}^\infty D_n$.

It remains to verify this works. Let $y^* \in X'$. I want to show that $\|y^*\| = \sup_{x \in D} |y^*(x)|$. There exists $y, \|y\| \leq 1$, such that

$$|y^*(y)| > \|y^*\| - \varepsilon.$$

By density, there exists one of the x_k^* from the countable dense subset of X' such that also

$$\|x_k^* - y^*\| < \varepsilon, \text{ so } |x_k^*(y)| > \|y^*\| - 2\varepsilon$$

Now $x_k^*(y) \in \phi_k(B)$ and so there exists $x \in D_k \subseteq D \subseteq B$ such that also

$$|x_k^*(x)| > \|y^*\| - 2\varepsilon.$$

Then since $\|x_k^* - y^*\| < \varepsilon$, this implies

$$\|y^*\| \geq |y^*(x)| = |(y^* - x_k^*)(x) + x_k^*(x)| \geq |x_k^*(x)| - \varepsilon > \|y^*\| - 3\varepsilon$$

It follows that

$$\|y^*\| - 3\varepsilon \leq \sup_{x \in D} |y^*(x)| \leq \|y^*\|$$

This proves the lemma because ε is arbitrary. ■

The next theorem is another version of the Pettis theorem. First here is a definition.

Definition 24.1.14 A function y having values in X' is weak * measurable, when for each $x \in X$, $y(\cdot)(x)$ is a measurable scalar valued function.

Theorem 24.1.15 If X' is separable and $y : \Omega \rightarrow X'$ is weak * measurable meaning $\omega \rightarrow y(\omega)(x)$ is a \mathbb{F} valued measurable function, then y is strongly measurable.

Proof: It is necessary to show $y^{-1}(B(a^*, r))$ is measurable for $a^* \in X'$. This will suffice because the separability of X' implies every open set is the countable union of such balls of the form $B(a^*, r)$. It also suffices to verify inverse images of closed balls are measurable because every open ball is the countable union of closed balls. From Lemma 24.1.13,

$$\begin{aligned} y^{-1}(\overline{B(a^*, r)}) &= \{\omega : \|y(\omega) - a^*\| \leq r\} \\ &= \left\{ \omega : \sup_{x \in D} |(y(\omega) - a^*)(x)| \leq r \right\} \\ &= \left\{ \omega : \sup_{x \in D} |y(\omega)(x) - a^*(x)| \leq r \right\} \\ &= \bigcap_{x \in D} y(\cdot)(x)^{-1}(\overline{B(a^*(x), r)}) \end{aligned}$$

which is a countable intersection of measurable sets by hypothesis. ■

The following are interesting consequences of the theory developed so far and are of interest independent of the theory of integration of vector valued functions.

Theorem 24.1.16 If X' is separable, then so is X .

Proof: Let $D = \{x_m\} \subseteq B$, the unit ball of X , be the sequence promised by Lemma 24.1.13. Let V be all finite linear combinations of elements of $\{x_m\}$ with rational scalars. Thus \bar{V} is a separable subspace of X . The claim is that $\bar{V} = X$. If not, then it follows that there exists $x_0 \in X \setminus \bar{V}$. But by the Hahn Banach theorem there exists $x_0^* \in X'$ satisfying $x_0^*(x_0) \neq 0$, but $x_0^*(v) = 0$ for every $v \in \bar{V}$. Hence $\|x_0^*\| = \sup_{x \in D} |x_0^*(x)| = 0$, a contradiction. ■

Corollary 24.1.17 If X is reflexive, then X is separable if and only if X' is separable.

Proof: From the above theorem, if X' is separable, then so is X . Now suppose X is separable with a dense subset equal to D . Then since X is reflexive, $J(D)$ is dense in X'' where J is the James map satisfying $Jx(x^*) \equiv x^*(x)$. Recall how this J preserves norms and maps onto X'' for X reflexive. Then since X'' is separable, it follows from the above theorem that X' is also separable. ■

Note how this shows that $L^1(\mathbb{R}^p, m_p)$ is not reflexive because this is a separable space, but $L^\infty(\mathbb{R}^p, m_p)$ is clearly not. For example, you could consider $\mathcal{X}_{[0, r]}$ for r a positive irrational number. There are uncountably many of these functions in $L^\infty([0, 1])$ and it is clear that $\|\mathcal{X}_{[0, r]} - \mathcal{X}_{[0, \bar{r}]\|_\infty = 1$.

24.1.1 Egoroff's Theorem

In the context of a more general notion of measurable function having values in a metric space, here is a version of Egoroff's theorem. Here we introduce a finite measure μ . None of the above section had anything to do with a measure.

Theorem 24.1.18 (Egoroff) *Let $(\Omega, \mathcal{F}, \mu)$ be a finite measure space, $(\mu(\Omega) < \infty)$ and let f_n, f be X valued measurable functions where X is a separable metric space and for all $\omega \notin E$ where $\mu(E) = 0, f_n(\omega) \rightarrow f(\omega)$. Then for every $\varepsilon > 0$, there exists a set, $F \supseteq E, \mu(F) < \varepsilon$, such that f_n converges uniformly to f on F^C .*

Proof: First suppose $E = \emptyset$ so that convergence is pointwise everywhere. Let

$$E_{km} = \{\omega \in \Omega : d(f_n(\omega), f(\omega)) \geq 1/m \text{ for some } n > k\}.$$

Claim: $[\omega : d(f_n(\omega), f(\omega)) \geq \frac{1}{m}]$ is measurable.

Proof of claim: Let $\{x_k\}_{k=1}^\infty$ be a countable dense subset of X and let r denote a positive rational number, \mathbb{Q}^+ . Then

$$\bigcup_{k \in \mathbb{N}, r \in \mathbb{Q}^+} f_n^{-1}(B(x_k, r)) \cap f^{-1}\left(B\left(x_k, \frac{1}{m} - r\right)\right) = \left[d(f, f_n) < \frac{1}{m}\right] \quad (24.2)$$

Here is why. If ω is in the set on the left, then $d(f_n(\omega), x_k) < r$ and $d(f(\omega), x_k) < \frac{1}{m} - r$. Therefore,

$$d(f(\omega), f_n(\omega)) < r + \frac{1}{m} - r = \frac{1}{m}.$$

Thus the left side is contained in the right. Now let ω be in the right side. That is $d(f_n(\omega), f(\omega)) < \frac{1}{m}$. Choose $2r < \frac{1}{m} - d(f_n(\omega), f(\omega))$ and pick $x_k \in B(f_n(\omega), r)$. Then

$$d(f(\omega), x_k) \leq d(f(\omega), f_n(\omega)) + d(f_n(\omega), x_k) < \frac{1}{m} - 2r + r = \frac{1}{m} - r$$

Thus $\omega \in f_n^{-1}(B(x_k, r)) \cap f^{-1}(B(x_k, \frac{1}{m} - r))$ and so ω is in the left side. Thus the two sets are equal. Now the set on the left in 24.2 is measurable because it is a countable union of measurable sets. This proves the claim since $[\omega : d(f_n(\omega), f(\omega)) \geq \frac{1}{m}]$ is the complement of this measurable set.

Hence E_{km} is measurable because $E_{km} = \bigcup_{n=k+1}^\infty [\omega : d(f_n(\omega), f(\omega)) \geq \frac{1}{m}]$. For fixed $m, \bigcap_{k=1}^\infty E_{km} = \emptyset$ because $f_n(\omega)$ converges to $f(\omega)$. Therefore, if $\omega \in \Omega$ there exists k such that if $n > k, |f_n(\omega) - f(\omega)| < \frac{1}{m}$ which means $\omega \notin E_{km}$. Note also that $E_{km} \supseteq E_{(k+1)m}$. Since $\mu(E_{1m}) < \infty$, Theorem 9.2.4 on Page 242 implies

$$0 = \mu(\bigcap_{k=1}^\infty E_{km}) = \lim_{k \rightarrow \infty} \mu(E_{km}).$$

Let $k(m)$ be chosen such that $\mu(E_{k(m)m}) < \varepsilon 2^{-m}$ and let $F = \bigcup_{m=1}^\infty E_{k(m)m}$. Then $\mu(F) < \varepsilon$ because

$$\mu(F) \leq \sum_{m=1}^\infty \mu(E_{k(m)m}) < \sum_{m=1}^\infty \varepsilon 2^{-m} = \varepsilon$$

Now let $\eta > 0$ be given and pick m_0 such that $m_0^{-1} < \eta$. If $\omega \in F^C$, then $\omega \in \bigcap_{m=1}^\infty E_{k(m)m}^C$.

Hence $\omega \in E_{k(m_0)m_0}^C$ so $d(f(\omega), f_n(\omega)) < 1/m_0 < \eta$ for all $n > k(m_0)$. This holds for all $\omega \in F^C$ and so f_n converges uniformly to f on F^C .

Now if $E \neq \emptyset$, consider $\{\mathcal{X}_{E^C} f_n\}_{n=1}^\infty$. Then $\mathcal{X}_{E^C} f_n$ is measurable and the sequence converges pointwise to $\mathcal{X}_E f$ everywhere. Therefore, from the first part, there exists a set of measure less than ε , F such that on F^C , $\{\mathcal{X}_{E^C} f_n\}$ converges uniformly to $\mathcal{X}_{E^C} f$. Therefore, on $(E \cup F)^C$, $\{f_n\}$ converges uniformly to f . ■

24.2 The Bochner Integral

24.2.1 Definition and Basic Properties

Definition 24.2.1 Let $a_k \in X$, a Banach space and let a simple function $\omega \rightarrow x(\omega)$ be

$$x(\omega) = \sum_{k=1}^n a_k \mathcal{X}_{E_k}(\omega) \quad (24.3)$$

where for each k , E_k is measurable and $\mu(E_k) < \infty$. Thus this is a measurable finite valued function zero off a set of finite measure. Then define

$$\int_{\Omega} x(\omega) d\mu \equiv \sum_{k=1}^n a_k \mu(E_k).$$

Proposition 24.2.2 Definition 24.2.1 is well defined, the integral is linear on simple functions and

$$\left\| \int_{\Omega} x(\omega) d\mu \right\| \leq \int_{\Omega} \|x(\omega)\| d\mu$$

whenever x is a simple function.

Proof: It suffices to verify that if $\sum_{k=1}^n a_k \mathcal{X}_{E_k}(\omega) = 0$, then $\sum_{k=1}^n a_k \mu(E_k) = 0$. Let $f \in X'$. Then

$$f\left(\sum_{k=1}^n a_k \mathcal{X}_{E_k}(\omega)\right) = \sum_{k=1}^n f(a_k) \mathcal{X}_{E_k}(\omega) = 0$$

and, therefore,

$$0 = \int_{\Omega} \left(\sum_{k=1}^n f(a_k) \mathcal{X}_{E_k}(\omega)\right) d\mu = \sum_{k=1}^n f(a_k) \mu(E_k) = f\left(\sum_{k=1}^n a_k \mu(E_k)\right).$$

Since $f \in X'$ is arbitrary, and X' separates the points of X , $\sum_{k=1}^n a_k \mu(E_k) = 0$ as hoped. It is now obvious that the integral is linear on simple functions.

As to the triangle inequality, say $x(\omega) = \sum_{k=1}^n a_k \mathcal{X}_{E_k}(\omega)$ where the E_k are disjoint. Then from the triangle inequality,

$$\left\| \int_{\Omega} x(\omega) d\mu \right\| = \left\| \sum_{k=1}^n a_k \mu(E_k) \right\| \leq \sum_{k=1}^n \|a_k\| \mu(E_k) = \int_{\Omega} \|x(\omega)\| d\mu \quad \blacksquare$$

Definition 24.2.3 A strongly measurable function x is Bochner integrable if there exists a sequence of simple functions x_n converging to x pointwise and satisfying

$$\int_{\Omega} \|x_n(\omega) - x_m(\omega)\| d\mu \rightarrow 0 \text{ as } m, n \rightarrow \infty. \quad (24.4)$$

If x is Bochner integrable, define

$$\int_{\Omega} x(\omega) d\mu \equiv \lim_{n \rightarrow \infty} \int_{\Omega} x_n(\omega) d\mu. \quad (24.5)$$

First it is important to show that this integral is well defined. When this is done, an easier to use condition will be developed. Note that by Lemma 24.1.2, if x is strongly measurable, $\|x\|$ is a measurable real valued function. Thus, it makes sense to consider $\int_{\Omega} \|x\| d\mu$ and also $\int_{\Omega} \|x - x_n\| d\mu$.

Theorem 24.2.4 *The definition of Bochner integrability is well defined. Also, a strongly measurable function x is Bochner integrable if and only if $\int_{\Omega} \|x\| d\mu < \infty$. In this case that the function is Bochner integrable, an approximating sequence of simple functions $\{y_n\}$ exists such that $\|y_n(\omega)\| \leq 2\|x(\omega)\|$ for all ω and*

$$\lim_{n \rightarrow \infty} \int_{\Omega} \|y_n(\omega) - x(\omega)\| d\mu = 0$$

Proof: \Rightarrow First consider the claim about the integral being well defined. Let $\{x_n\}$ be a sequence of simple functions converging pointwise to x and satisfying the conditions given above for x to be Bochner integrable. Then

$$\left| \int_{\Omega} \|x_n(\omega)\| d\mu - \int_{\Omega} \|x_m(\omega)\| d\mu \right| \leq \int_{\Omega} \|x_n - x_m\| d\mu$$

which is given to converge to 0 as $n, m \rightarrow \infty$ which shows that $\{\int_{\Omega} \|x_n(\omega)\| d\mu\}_{n=1}^{\infty}$ is a Cauchy sequence. Hence it is bounded and so, by Fatou's lemma,

$$\int_{\Omega} \|x(\omega)\| d\mu \leq \liminf_{n \rightarrow \infty} \int_{\Omega} \|x_n(\omega)\| d\mu < \infty$$

The limit in 24.5 exists because

$$\left\| \int_{\Omega} x_n d\mu - \int_{\Omega} x_m d\mu \right\| = \left\| \int_{\Omega} (x_n - x_m) d\mu \right\| \leq \int_{\Omega} \|x_n - x_m\| d\mu$$

and the last term is no more than ε whenever n, m are large enough. From Fatou's lemma, if n is large enough,

$$\int_{\Omega} \|x_n - x\| d\mu < \varepsilon$$

Now if you have another sequence $\{\hat{x}_n\}$ satisfying the condition 24.4 along with pointwise convergence to x ,

$$\begin{aligned} \left\| \int_{\Omega} x_n d\mu - \int_{\Omega} \hat{x}_n d\mu \right\| &= \left\| \int_{\Omega} (x_n - \hat{x}_n) d\mu \right\| \leq \int_{\Omega} \|x_n - \hat{x}_n\| d\mu \\ &\leq \int_{\Omega} \|x_n - x\| d\mu + \int_{\Omega} \|x - \hat{x}_n\| d\mu < 2\varepsilon \end{aligned}$$

if n is large enough. Hence convergence of the integrals of the simple functions takes place and these integrals converge to the same thing. Thus the definition is well defined and $\int_{\Omega} \|x\| d\mu < \infty$.

\Leftarrow Next suppose $\int_{\Omega} \|x\| d\mu < \infty$ for x strongly measurable. By Lemma 24.1.2, there is a sequence of finite valued measurable functions $\{y_n\}$ with $\|y_n(\omega)\| \leq 2\|x(\omega)\|$ and

$y_n(\omega) \rightarrow x(\omega)$ for each ω . Thus, in fact, y_n is a simple function because it must be zero off a set of finite measure because

$$\int_{\Omega} \|y_n(\omega)\| d\mu < 2 \int_{\Omega} \|x(\omega)\| d\mu$$

Then by the dominated convergence theorem for scalar valued functions,

$$\lim_{n \rightarrow \infty} \int_{\Omega} \|y_n - x\| d\mu = 0$$

Thus,

$$\int_{\Omega} \|y_n - y_m\| d\mu \leq \int_{\Omega} \|y_n - x\| d\mu + \int_{\Omega} \|x - y_m\| d\mu < \varepsilon$$

if m, n are large enough so $\{y_n\}$ is a suitable approximating sequence for x . ■

This is a very nice theorem. It says that all you have to do is verify measurability and absolute integrability just like the case of scalar valued functions. Other things which are totally similar are that the integral is linear, the triangle inequality holds, and you can take a continuous linear functional inside the integral. These things are considered in the following theorem.

Theorem 24.2.5 *The Bochner integral is well defined and if x is Bochner integrable and $f \in X'$,*

$$f\left(\int_{\Omega} x(\omega) d\mu\right) = \int_{\Omega} f(x(\omega)) d\mu \quad (24.6)$$

and the triangle inequality is valid,

$$\left\|\int_{\Omega} x(\omega) d\mu\right\| \leq \int_{\Omega} \|x(\omega)\| d\mu. \quad (24.7)$$

Also, the Bochner integral is linear. That is, if a, b are scalars and x, y are two Bochner integrable functions, then

$$\int_{\Omega} (ax(\omega) + by(\omega)) d\mu = a \int_{\Omega} x(\omega) d\mu + b \int_{\Omega} y(\omega) d\mu \quad (24.8)$$

Proof: Theorem 24.2.4 shows $\int_{\Omega} \|x(\omega)\| d\mu < \infty$ and that the definition of the integral is well defined.

It remains to verify the triangle inequality on Bochner integral functions and the claim about passing a continuous linear functional inside the integral. First of all, consider the triangle inequality. From Lemma 24.1.2, there is a sequence of simple functions $\{y_n\}$ satisfying 24.4 and converging to x pointwise such that also $\|y_n(\omega)\| \leq 2\|x(\omega)\|$. Thus,

$$\left\|\int_{\Omega} x(\omega) d\mu\right\| \equiv \lim_{n \rightarrow \infty} \left\|\int_{\Omega} y_n(\omega) d\mu\right\| \leq \lim_{n \rightarrow \infty} \int_{\Omega} \|y_n(\omega)\| d\mu = \int_{\Omega} \|x(\omega)\| d\mu$$

the last step coming from the dominated convergence theorem since $\|y_n(\omega)\| \leq 2\|x(\omega)\|$ and $\|y_n(\omega)\| \rightarrow \|x(\omega)\|$ for each ω . This shows the triangle inequality.

From Definition 24.2.1 and Theorem 24.2.4 and $\{y_n\}$ being the approximating sequence described there,

$$f\left(\int_{\Omega} y_n d\mu\right) = \int_{\Omega} f(y_n) d\mu.$$

Thus,

$$f\left(\int_{\Omega} x d\mu\right) = \lim_{n \rightarrow \infty} f\left(\int_{\Omega} y_n d\mu\right) = \lim_{n \rightarrow \infty} \int_{\Omega} f(y_n) d\mu = \int_{\Omega} f(x) d\mu,$$

the last equation holding from the dominated convergence theorem ($|f(y_n)| \leq \|f\| \|y_n\| \leq 2\|f\| \|x\|$). This shows 24.6.

It remains to verify 24.8. Let $f \in X'$. Then from 24.6

$$\begin{aligned} f\left(\int_{\Omega} (ax(\omega) + by(\omega)) d\mu\right) &= \int_{\Omega} (af(x(\omega)) + bf(y(\omega))) d\mu \\ &= a \int_{\Omega} f(x(\omega)) d\mu + b \int_{\Omega} f(y(\omega)) d\mu \\ &= f\left(a \int_{\Omega} x(\omega) d\mu + b \int_{\Omega} y(\omega) d\mu\right). \end{aligned}$$

Since X' separates the points of X , it follows

$$\int_{\Omega} (ax(\omega) + by(\omega)) d\mu = a \int_{\Omega} x(\omega) d\mu + b \int_{\Omega} y(\omega) d\mu$$

and this proves 24.8. ■

A similar result is the following corollary.

Corollary 24.2.6 *Let an X valued function x be Bochner integrable. Let $L \in \mathcal{L}(X, Y)$ where Y is another Banach space. Then Lx is a Y valued Bochner integrable function and*

$$L\left(\int_{\Omega} x(\omega) d\mu\right) = \int_{\Omega} Lx(\omega) d\mu$$

Proof: From Theorem 24.2.4 there is a sequence of simple functions $\{y_n\}$ having the properties listed in that theorem. These are measurable with finitely many values and are forced to be simple because $\|y_n\| \leq 2\|x\|$. Then consider $\{Ly_n\}$ which converges pointwise to Lx . Since L is continuous and linear,

$$\int_{\Omega} \|Ly_n - Lx\|_Y d\mu \leq \|L\| \int_{\Omega} \|y_n - x\|_X d\mu$$

which converges to 0. This implies

$$\lim_{m, n \rightarrow \infty} \int_{\Omega} \|Ly_n - Ly_m\| d\mu = 0$$

and so by definition Lx is Bochner integrable. Also

$$\begin{aligned} \int_{\Omega} x(\omega) d\mu &= \lim_{n \rightarrow \infty} \int_{\Omega} y_n(\omega) d\mu \\ \int_{\Omega} Lx(\omega) d\mu &= \lim_{n \rightarrow \infty} \int_{\Omega} Ly_n(\omega) d\mu = \lim_{n \rightarrow \infty} L \int_{\Omega} y_n(\omega) d\mu \end{aligned}$$

Next,

$$\begin{aligned} &\left\| L\left(\int_{\Omega} x(\omega) d\mu\right) - \int_{\Omega} Lx(\omega) d\mu \right\|_Y \\ &\leq \left\| L\left(\int_{\Omega} x(\omega) d\mu\right) - L \int_{\Omega} y_n(\omega) d\mu \right\|_Y \\ &+ \left\| \int_{\Omega} Ly_n(\omega) d\mu - \int_{\Omega} Lx(\omega) d\mu \right\|_Y < \varepsilon/2 + \varepsilon/2 = \varepsilon \end{aligned}$$

whenever n large enough. ■

24.2.2 Taking a Closed Operator Out of the Integral

Now let X and Y be separable Banach spaces and suppose $A : D(A) \subseteq X \rightarrow Y$ be a closed operator. Recall this means that the graph of A ,

$$G(A) \equiv \{(x, Ax) : x \in D(A)\}$$

is a closed subset of $X \times Y$ with respect to the product topology obtained from the norm

$$\|(x, y)\| = \max(\|x\|, \|y\|).$$

Thus also $G(A)$ is a separable Banach space with the above norm. You can also consider $D(A)$ as a separable Banach space having the graph norm

$$\|x\|_{D(A)} \equiv \max(\|x\|, \|Ax\|) \quad (24.9)$$

which is isometric to $G(A)$ with the mapping, $\theta x \equiv (x, Ax)$. Recall why this is. It is clear that θ is one to one and onto $G(A)$. Is it continuous? If $x_n \rightarrow x$ in $D(A)$, this means that $x_n \rightarrow x$ in X and $Ax_n \rightarrow y$. Then, since A is closed, it follows that $y = Ax$ so $(x_n, Ax_n) \rightarrow (x, Ax)$ in $G(A)$. Hence θ is indeed continuous and onto. Similar reasoning shows that $D(A)$ with this norm is complete. Hence it is a Banach space. Thus θ^{-1} is also continuous. The following lemma is a fundamental result which was proved earlier in the discussion on the Eberlein Smulian theorem in which this was an essential fact to allow the case of a reflexive Banach space which maybe was not separable. See Lemma 21.5.13 for the proof.

Lemma 24.2.7 *A closed subspace of a reflexive Banach space is reflexive.*

Then, with this lemma, one has the following corollary.

Corollary 24.2.8 *Suppose Y is a reflexive Banach space and X is a Banach space such that there exists a continuous one to one mapping, $g : X \rightarrow Y$ such that $g(X)$ is a closed subset of Y . Then X is reflexive.*

Proof: By the open mapping theorem, $g(X)$ and X are homeomorphic since g^{-1} must also be continuous. Therefore, since $g(X)$ is reflexive because it is a closed subspace of a reflexive space, it follows X is also reflexive. ■

Lemma 24.2.9 *Suppose V is a reflexive Banach space and that V is a dense subset of W , another Banach space in the topology of W . Then i^*W' is a dense subset of V' where here i is the inclusion map of V into W .*

Proof: First note that i^* is one to one. If $i^*w^* = 0$ for $w^* \in W'$, then this means that for all $v \in V$,

$$i^*w(v) = w^*(v) = 0$$

and since V is dense in W , this shows $w^* = 0$.

Consider the following diagram

$$\begin{array}{ccc} V'' & \xrightarrow{i^{**}} & W'' \\ V' & \xleftarrow{i^*} & W' \\ V & \xrightarrow{i} & W \end{array}$$

in which i is the inclusion map. Next suppose i^*W' is not dense in V' . Then, using the Hahn Banach theorem, there exists $v^{**} \in V''$ such that $v^{**} \neq 0$ but $v^{**}(i^*W') = 0$. It follows from V being reflexive, that $v^{**} = Jv_0$ where J is the James map from V to V'' for some $v_0 \in V$. Thus for every $w^* \in W'$,

$$\begin{aligned} 0 &= v^{**}(i^*w^*) \equiv i^{**}v^{**}(w^*) \\ &= i^{**}Jv_0(w^*) = Jv_0(i^*w^*) \\ &\equiv i^*w^*(v_0) = w^*(v_0) \end{aligned}$$

and since W' separates the points of W , it follows $v_0 = 0$ which contradicts $v^{**} \neq 0$. ■

Note that in the proof, only V reflexive was used.

This lemma implies an easy corollary.

Corollary 24.2.10 *Let E and F be reflexive Banach spaces and let A be a closed operator $A : D(A) \subseteq E \rightarrow F$. Suppose also that $D(A)$ is dense in E . Then making $D(A)$ into a Banach space by using the above graph norm given in 24.9, it follows that $D(A)$ is a Banach space and i^*E' is a dense subspace of $D(A)'$.*

Proof: First note that $E \times F$ is a reflexive Banach space and $\mathcal{G}(A)$ is a closed subspace of $E \times F$ so it is also a reflexive Banach space. Now $D(A)$ is isometric to $\mathcal{G}(A)$ and so it follows $D(A)$ is a dense subspace of E which is reflexive. Therefore, from Lemma 24.2.9 the conclusion follows. ■

With this preparation, here is another interesting theorem. This one is about taking outside the integral a closed linear operator as opposed to a continuous linear operator.

Theorem 24.2.11 *Let X, Y be separable Banach spaces and let $A : D(A) \subseteq X \rightarrow Y$ be a closed operator where $D(A)$ is a dense separable subset of X with respect to the graph norm on $D(A)$ described above¹. Suppose also that i^*X' is a dense subspace of $D(A)'$ where $D(A)$ is a Banach space having the graph norm described in 24.9. Suppose that $(\Omega, \mathcal{F}, \mu)$ is a measure space and $x : \Omega \rightarrow X$ is strongly measurable and it happens that $x(\omega) \in D(A)$ for all $\omega \in \Omega$. Then x is strongly measurable as a mapping into $D(A)$. Also Ax is strongly measurable as a map into Y and if*

$$\int_{\Omega} \|x(\omega)\| d\mu, \int_{\Omega} \|Ax(\omega)\| d\mu < \infty, \quad (24.10)$$

then

$$\int_{\Omega} x(\omega) d\mu \in D(A) \quad (24.11)$$

and

$$A \int_{\Omega} x(\omega) d\mu = \int_{\Omega} Ax(\omega) d\mu. \quad (24.12)$$

Proof: First of all, consider the assertion that x is strongly measurable into $D(A)$. Letting $f \in D(A)'$ be given, there exists a sequence, $\{g_n\} \subseteq i^*X'$ such that $g_n \rightarrow f$ in $D(A)'$. Therefore, $\omega \rightarrow g_n(x(\omega))$ is measurable by assumption and $g_n(x(\omega)) \rightarrow f(x(\omega))$, which shows that $\omega \rightarrow f(x(\omega))$ is measurable. By the Pettis theorem, it follows that $\omega \rightarrow x(\omega)$ is strongly measurable as a map into $D(A)$.

¹Note that this follows from the assumed separability of X, Y because the graph is a subset of the separable space $X \times Y$

It follows from Theorem 24.2.4 there exists a sequence of simple functions $\{x_n\}$ of the form

$$x_n(\omega) = \sum_{k=1}^{m_n} a_k^n \mathcal{K}_{E_k^n}(\omega), x_n(\omega) \in D(A),$$

which converges strongly and pointwise to $x(\omega)$ in $D(A)$. Thus

$$x_n(\omega) \rightarrow x(\omega), Ax_n(\omega) \rightarrow Ax(\omega),$$

which shows $\omega \rightarrow Ax(\omega)$ is strongly measurable in Y as claimed.

It remains to verify the assertions about the integral. 24.10 implies x is Bochner integrable as a function having values in $D(A)$ with the norm on $D(A)$ described above. Therefore, by Theorem 24.2.4 there exists a sequence of simple functions $\{y_n\}$ having values in $D(A)$, $\lim_{m,n \rightarrow \infty} \int_{\Omega} \|y_n - y_m\|_{D(A)} d\mu = 0$, $y_n(\omega)$ converging pointwise to $x(\omega)$, and also $\|y_n(\omega)\|_{D(A)} \leq 2\|x(\omega)\|_{D(A)}$ and $\lim_{n \rightarrow \infty} \int_{\Omega} \|x(\omega) - y_n(\omega)\|_{D(A)} ds = 0$. Therefore,

$$\int_{\Omega} y_n(\omega) d\mu \in D(A), \int_{\Omega} y_n(\omega) d\mu \rightarrow \int_{\Omega} x(\omega) d\mu \text{ in } X,$$

and since y_n is a simple function and A is linear,

$$A \int_{\Omega} y_n(\omega) d\mu = \int_{\Omega} Ay_n(\omega) d\mu \rightarrow \int_{\Omega} Ax(\omega) d\mu \text{ in } Y.$$

It follows, since A is a closed operator, that $\int_{\Omega} x(\omega) d\mu \in D(A)$ and

$$A \int_{\Omega} x(\omega) d\mu = \int_{\Omega} Ax(\omega) d\mu. \blacksquare$$

Here is another version of this theorem which has different hypotheses.

Theorem 24.2.12 *Let X and Y be separable Banach spaces and let $A : D(A) \subseteq X \rightarrow Y$ be a closed operator. Also let $(\Omega, \mathcal{F}, \mu)$ be a measure space and let $x : \Omega \rightarrow X$ be Bochner integrable such that $x(\omega) \in D(A)$ for all ω . Also suppose Ax is Bochner integrable. Then*

$$\int_{\Omega} Ax d\mu = A \int_{\Omega} x d\mu$$

and $\int_{\Omega} x d\mu \in D(A)$.

Proof: Consider the graph of A ,

$$G(A) \equiv \{(x, Ax) : x \in D(A)\} \subseteq X \times Y.$$

Then since A is closed, $G(A)$ is a closed separable Banach space with the norm $\|(x, y)\| \equiv \max(\|x\|, \|y\|)$. Therefore, for $g^* \in G(A)'$, apply the Hahn Banach theorem and obtain $(x^*, y^*) \in (X \times Y)'$ such that $g^*(x, Ax) = (x^*(x), y^*(Ax))$. Now it follows from the assumptions that $\omega \rightarrow (x^*(x(\omega)), y^*(Ax(\omega)))$ is measurable with values in $G(A)$. It is also separably valued because this is true of $G(A)$. By the Pettis theorem, $\omega \rightarrow (x(\omega), Ax(\omega))$ must be strongly measurable. Also $\int \|x(\omega)\| + \|Ax(\omega)\| d\mu < \infty$ by assumption and so there exists a sequence of simple functions having values in $G(A)$, $\{(x_n(\omega), Ax_n(\omega))\}$

which converges to $(x(\omega), A(\omega))$ pointwise such that $\int \|(x_n, Ax_n) - (x, Ax)\| d\mu \rightarrow 0$ in $G(A)$. Now for simple functions it is routine to verify that

$$\int (x_n, Ax_n) d\mu = \left(\int x_n d\mu, \int Ax_n d\mu \right) = \left(\int x_n d\mu, A \int x_n d\mu \right)$$

Also

$$\begin{aligned} \left\| \int x_n d\mu - \int x d\mu \right\| &\leq \int \|x_n - x\| d\mu \\ &\leq \int \|(x_n, Ax_n) - (x, Ax)\| d\mu \end{aligned}$$

which converges to 0. Also

$$\begin{aligned} \left\| \int Ax_n d\mu - \int Ax d\mu \right\| &= \left\| A \int x_n d\mu - A \int x d\mu \right\| \\ &\leq \int \|Ax_n - Ax\| d\mu \\ &\leq \int \|(x_n, Ax_n) - (x, Ax)\| d\mu \end{aligned}$$

and this converges to 0. Therefore, $\int x_n d\mu \rightarrow \int x d\mu$ and $A \int x_n d\mu \rightarrow \int Ax d\mu$. Since each $\int x_n d\mu \in D(A)$, and A is closed, this implies $\int x d\mu \in D(A)$ and $A \int x d\mu = \int Ax d\mu$. ■

24.3 Operator Valued Functions

Consider the case where $A(\omega) \in \mathcal{L}(X, Y)$ for X and Y separable Banach spaces. With the operator norm $\mathcal{L}(X, Y)$ is a Banach space and so if A is strongly measurable, the Bochner integral can be defined as before. However, it is also possible to define the Bochner integral of such operator valued functions for more general situations. In this section, $(\Omega, \mathcal{F}, \mu)$ will be a measure space as usual.

Lemma 24.3.1 *Let $x \in X$ and suppose A is strongly measurable. Then $\omega \rightarrow A(\omega)x$ is strongly measurable as a map into Y .*

Proof: Since A is assumed to be strongly measurable, it is the pointwise limit of measurable finite valued functions of the form $A_n(\omega) \equiv \sum_{k=1}^{m_n} A_k^n \chi_{E_k^n}(\omega)$ where A_k^n is in $\mathcal{L}(X, Y)$. It follows $A_n(\omega)x \rightarrow A(\omega)x$ for each ω and so, since $\omega \rightarrow A_n(\omega)x$ is a simple Y valued function, $\omega \rightarrow A(\omega)x$ must be strongly measurable. ■

Definition 24.3.2 *Suppose $A(\omega) \in \mathcal{L}(X, Y)$ for each $\omega \in \Omega$ where X, Y are separable Banach spaces. Suppose also that for each $x \in X$,*

$$\omega \rightarrow A(\omega)x \text{ is strongly measurable} \quad (24.13)$$

and there exists C such that for each $x \in X$,

$$\int_{\Omega} \|A(\omega)x\| d\mu < C\|x\| \quad (24.14)$$

Then $\int_{\Omega} A(\omega) d\mu \in \mathcal{L}(X, Y)$ is defined by the following formula.

$$\left(\int_{\Omega} A(\omega) d\mu \right)(x) \equiv \int_{\Omega} A(\omega)x d\mu \quad (24.15)$$

Lemma 24.3.3 *The above definition is well defined. Furthermore, if 24.13 holds then $\omega \rightarrow \|A(\omega)\|$ is measurable and if 24.14 holds, then $\|\int_{\Omega} A(\omega) d\mu\| \leq \int_{\Omega} \|A(\omega)\| d\mu$.*

Proof: It is clear that in case $\omega \rightarrow A(\omega)x$ is measurable for all $x \in X$ there exists a unique $\Psi \in \mathcal{L}(X, Y)$ such that $\Psi(x) = \int_{\Omega} A(\omega)x d\mu$. This is because $x \rightarrow \int_{\Omega} A(\omega)x d\mu$ is linear and continuous. It is continuous because

$$\left\| \int_{\Omega} A(\omega)x d\mu \right\| \leq \int_{\Omega} \|A(\omega)x\| d\mu \leq \int_{\Omega} \|A(\omega)\| d\mu \|x\|$$

Thus $\Psi = \int_{\Omega} A(\omega) d\mu$ and the definition is well defined.

Now consider the assertion about $\omega \rightarrow \|A(\omega)\|$. Let $D' \subseteq B'$ the closed unit ball in Y' be such that D' is countable and $\|y\| = \sup_{y^* \in D'} |y^*(y)|$. This is from Lemma 24.1.7. Recall X is separable. Also let D be a countable dense subset of B , the unit ball of X . Then

$$\begin{aligned} \{\omega : \|A(\omega)\| > \alpha\} &= \left\{ \omega : \sup_{x \in D} \|A(\omega)x\| > \alpha \right\} = \cup_{x \in D} \{\omega : \|A(\omega)x\| > \alpha\} \\ &= \cup_{x \in D} (\cup_{y^* \in D'} \{|y^*(A(\omega)x)| > \alpha\}) \end{aligned}$$

and this is measurable because $\omega \rightarrow A(\omega)x$ is strongly, hence weakly measurable.

Now suppose 24.14 holds. Then for all x , $\int_{\Omega} \|A(\omega)x\| d\mu < C\|x\|$. It follows that for $\|x\| \leq 1$,

$$\left\| \left(\int_{\Omega} A(\omega) d\mu \right) (x) \right\| = \left\| \int_{\Omega} A(\omega)x d\mu \right\| \leq \int_{\Omega} \|A(\omega)x\| d\mu \leq \int_{\Omega} \|A(\omega)\| d\mu$$

and so $\|\int_{\Omega} A(\omega) d\mu\| \leq \int_{\Omega} \|A(\omega)\| d\mu$. ■

Now it is interesting to consider the case where $A(\omega) \in \mathcal{L}(H, H)$ where $\omega \rightarrow A(\omega)x$ is strongly measurable and $A(\omega)$ is compact and self adjoint. Recall the Kuratowski measurable selection theorem, Theorem 9.15.8 on Page 274 listed here for convenience.

Theorem 24.3.4 *Let E be a compact metric space and let (Ω, \mathcal{F}) be a measure space. Suppose $\psi : E \times \Omega \rightarrow \mathbb{R}$ has the property that $x \rightarrow \psi(x, \omega)$ is continuous and $\omega \rightarrow \psi(x, \omega)$ is measurable. Then there exists a measurable function, f having values in E such that $\psi(f(\omega), \omega) = \sup_{x \in E} \psi(x, \omega)$. Furthermore, $\omega \rightarrow \psi(f(\omega), \omega)$ is measurable.*

24.3.1 Review of Hilbert Schmidt Theorem

This section is a review of earlier material and is presented a little differently. I think it does not hurt to repeat some things relative to Hilbert space. I will give a proof of the Hilbert Schmidt theorem which will generalize to a result about measurable operators. It will be a little different then the earlier proof. Recall the following.

Definition 24.3.5 *Define $v \otimes u \in \mathcal{L}(H, H)$ by $v \otimes u(x) = (x, u)v$. $A \in \mathcal{L}(H, H)$ is a compact operator if whenever $\{x_k\}$ is a bounded sequence, there exists a convergent subsequence of $\{Ax_k\}$. Equivalently, A maps bounded sets to sets whose closures are compact or to use other terminology, A maps bounded sets to sets which are precompact.*

Next is a convenient description of compact operators on a Hilbert space.

Lemma 24.3.6 *Let H be a Hilbert space and suppose $A \in \mathcal{L}(H, H)$ is a compact operator. Then*

1. *A is a compact operator if and only if whenever $x_n \rightarrow x$ weakly in H , it follows that $Ax_n \rightarrow Ax$ strongly in H .*
2. *For $u, v \in H$, $v \otimes u : H \rightarrow H$ is a compact operator.*
3. *Let B be the closed unit ball in H . If A is self adjoint and compact, then if $x_n \rightarrow x$ weakly on B , it follows that $(Ax_n, x_n) \rightarrow (Ax, x)$ so $x \rightarrow |(Ax, x)|$ achieves its maximum value on B .*
4. *The function, $v \otimes u$ is compact and the operator $u \otimes u$ is self adjoint.*

Proof: Consider \Rightarrow of 1. Suppose then that $x_n \rightarrow x$ weakly. Since $\{x_n\}$ is weakly bounded, it follows from the uniform boundedness principle that $\{\|x_n\|\}$ is bounded. Let $x_n \in \hat{B}$ for \hat{B} some closed ball. If Ax_n fails to converge to Ax , then there is $\varepsilon > 0$ and a subsequence still denoted as $\{x_n\}$ such that $x_n \rightarrow x$ weakly but $\|Ax_n - Ax\| \geq \varepsilon > 0$. Then $A(\hat{B})$ is precompact because A is compact so there is a further subsequence, still denoted by $\{x_n\}$ such that Ax_n converges to some $y \in H$. Therefore,

$$\begin{aligned} (y, w) &= \lim_{n \rightarrow \infty} (Ax_n, w) = \lim_{n \rightarrow \infty} (x_n, A^*w) \\ &= (x, A^*w) = (Ax, w) \end{aligned}$$

which shows $Ax = y$ since w is arbitrary. However, this is a contradiction to $\|Ax_n - Ax\| \geq \varepsilon > 0$.

Consider \Leftarrow of 1. Why is A compact if it satisfies the property that it takes weakly convergent sequences to strongly convergent ones? If A is not compact, then there exists \hat{B} a bounded set such that $A(\hat{B})$ is not precompact. Thus, there exists a sequence $\{Ax_n\}_{n=1}^\infty \subseteq A(\hat{B})$ which has no convergent subsequence where $x_n \in \hat{B}$ the bounded set. However, there is a subsequence $\{x_n\} \in \hat{B}$ which converges weakly to some $x \in H$ because of weak compactness. Hence $Ax_n \rightarrow Ax$ by assumption and so this is a contradiction to there being no convergent subsequence of $\{Ax_n\}_{n=1}^\infty$.

Next consider 2. Letting $\{x_n\}$ be a bounded sequence,

$$v \otimes u(x_n) = (x_n, u)v.$$

There exists a weakly convergent subsequence of $\{x_n\}$ say $\{x_{n_k}\}$ converging weakly to $x \in H$. Therefore,

$$\|v \otimes u(x_{n_k}) - v \otimes u(x)\| = \|(x_{n_k}, u) - (x, u)\| \|v\|$$

which converges to 0. Thus $v \otimes u$ is compact as claimed. It takes bounded sets to precompact sets.

Next consider 3. To verify the assertion about $x \rightarrow (Ax, x)$, let $x_n \rightarrow x$ weakly. Since A is compact, $Ax_n \rightarrow Ax$ by part 1. Then, since A is self adjoint,

$$\begin{aligned} & |(Ax_n, x_n) - (Ax, x)| \\ \leq & |(Ax_n, x_n) - (Ax, x_n)| + |(Ax, x_n) - (Ax, x)| \\ \leq & |(Ax_n, x_n) - (Ax, x_n)| + |(Ax_n, x) - (Ax, x)| \\ \leq & \|Ax_n - Ax\| \|x_n\| + \|Ax_n - Ax\| \|x\| \leq 2 \|Ax_n - Ax\| \end{aligned}$$

which converges to 0. Now let $\{x_n\}$ be a maximizing sequence for $|(Ax, x)|$ for $x \in B$ and let $\lambda \equiv \sup \{|(Ax, x)| : x \in B\}$. There is a subsequence still denoted as $\{x_n\}$ which converges weakly to some $x \in B$ by weak compactness. Hence $|(Ax, x)| = \lim_{n \rightarrow \infty} |(Ax_n, x_n)| = \lambda$.

Next consider 4. It only remains to verify that $u \otimes u$ is self adjoint. This follows from the definition.

$$\begin{aligned} ((u \otimes u)x, y) &\equiv (u(x, u), y) = (x, u)(u, y) \\ (x, (u \otimes u)y) &\equiv (x, u(y, u)) = (u, y)(x, u), \end{aligned}$$

the same thing. ■

Observation 24.3.7 *Note that if A is any self adjoint operator,*

$$\overline{(Ax, x)} = (x, Ax) = (Ax, x).$$

so (Ax, x) is real valued.

From Lemma 24.3.6, the maximum of $|(Ax, x)|$ exists on the closed unit ball B .

Lemma 24.3.8 *Let $A \in \mathcal{L}(H, H)$ and suppose it is self adjoint and compact. Let B denote the closed unit ball in H . Let $e \in B$ be such that $|(Ae, e)| = \max_{x \in B} |(Ax, x)|$. Then letting $\lambda = (Ae, e)$, it follows $Ae = \lambda e$. You can always assume $\|e\| = 1$.*

Proof: From the above observation, (Ax, x) is always real and since A is compact, $|(Ax, x)|$ achieves a maximum at e . It remains to verify e is an eigenvector. If $|(Ae, e)| = 0$ for all $e \in B$, then A is a self adjoint nonnegative ($(Ax, x) \geq 0$) operator and so by Cauchy Schwarz inequality, $(Ae, x) \leq (Ax, x)^{1/2} (Ae, e)^{1/2} = 0$ and so $Ae = 0$ for all e . Assume then that A is not 0. You can always make $|(Ae, e)|$ at least as large by replacing e with $e/\|e\|$. Thus, there is no loss of generality in letting $\|e\| = 1$ in every case.

Suppose $\lambda = (Ae, e) \geq 0$ where $|(Ae, e)| = \max_{x \in B} |(Ax, x)|$. Thus

$$((\lambda I - A)e, e) = \lambda \|e\|^2 - \lambda = 0$$

Then it is easy to verify that $\lambda I - A$ is a nonnegative ($((\lambda I - A)x, x) \geq 0$ for all x .) and self adjoint operator. To see this, note that

$$((\lambda I - A)x, x) = \|x\|^2 \left((\lambda I - A) \frac{x}{\|x\|}, \frac{x}{\|x\|} \right) = \|x\|^2 \lambda - \|x\|^2 \left(A \frac{x}{\|x\|}, \frac{x}{\|x\|} \right) \geq 0$$

Therefore, the Cauchy Schwarz inequality can be applied to write

$$((\lambda I - A)e, x) \leq ((\lambda I - A)e, e)^{1/2} ((\lambda I - A)x, x)^{1/2} = 0$$

Since this is true for all x it follows $Ae = \lambda e$. Just pick $x = (\lambda I - A)e$.

Next suppose $\max_{x \in B} |(Ax, x)| = -(Ae, e)$. Let $-\lambda = -(Ae, e)$ and the previous result can be applied to $-A$ and $-\lambda$. Thus $-\lambda e = -Ae$ and so $Ae = \lambda e$. ■

With these lemmas here is a major theorem, the Hilbert Schmidt theorem. I think this proof is a little slicker than the more standard proof given earlier.

Theorem 24.3.9 *Let $A \in \mathcal{L}(H, H)$ be a compact self adjoint operator on a Hilbert space. Then there exist real numbers $\{\lambda_k\}_{k=1}^\infty$ and vectors $\{e_k\}_{k=1}^\infty$ such that*

$$\begin{aligned} \|e_k\| &= 1, (e_k, e_j)_H = 0 \text{ if } k \neq j, Ae_k = \lambda_k e_k, \\ |\lambda_n| &\geq |\lambda_{n+1}| \text{ for all } n, \lim_{n \rightarrow \infty} \lambda_n = 0, \\ \lim_{n \rightarrow \infty} \left\| A - \sum_{k=1}^n \lambda_k (e_k \otimes e_k) \right\|_{\mathcal{L}(H, H)} &= 0. \end{aligned} \quad (24.16)$$

Proof: This is done by considering a sequence of compact self adjoint operators, A, A_1, A_2, \dots . Here is how these are defined. Using Lemma 24.3.8 let e_1, λ_1 be given by that lemma such that

$$|(Ae_1, e_1)| = \max_{x \in B} |(Ax, x)|, \lambda_1 = (Ae_1, e_1) \Rightarrow Ae_1 = \lambda_1 e_1$$

Then by that lemma, $Ae_1 = \lambda_1 e_1$ and $\|e_1\| = 1$. Now define $A_1 = A - \lambda_1 e_1 \otimes e_1$. This is compact and self adjoint by Lemma 24.3.6. Thus, one could repeat the argument.

If A_n has been obtained, use Lemma 24.3.8 to obtain e_{n+1} and λ_{n+1} such that

$$|(A_n e_{n+1}, e_{n+1})| = \max_{x \in B} |(A_n x, x)|, \lambda_{n+1} = (A_n e_{n+1}, e_{n+1}).$$

By that lemma again, $A_n e_{n+1} = \lambda_{n+1} e_{n+1}$ and $\|e_{n+1}\| = 1$. Then $A_{n+1} \equiv A_n - \lambda_{n+1} e_{n+1} \otimes e_{n+1}$. Thus iterating this,

$$A_n = A - \sum_{k=1}^n \lambda_k e_k \otimes e_k. \quad (24.17)$$

Assume for $j, k \leq n, (e_k, e_j) = \delta_{jk}$. Then the new vector e_{n+1} will be orthogonal to the earlier ones. This is the next claim.

Claim 1: If $k < n+1$ then $(e_{n+1}, e_k) = 0$. Also $Ae_k = \lambda_k e_k$ for all k and from the construction, $A_n e_{n+1} = \lambda_{n+1} e_{n+1}$.

Proof of claim: From the above,

$$\lambda_{n+1} e_{n+1} = A_n e_{n+1} = Ae_{n+1} - \sum_{k=1}^n \lambda_k (e_{n+1}, e_k) e_k.$$

From the above and induction hypothesis that $(e_k, e_j) = \delta_{jk}$ for $j, k \leq n$,

$$\begin{aligned} \lambda_{n+1} (e_{n+1}, e_j) &= (Ae_{n+1}, e_j) - \sum_{k=1}^n \lambda_k (e_{n+1}, e_k) (e_k, e_j) \\ &= (e_{n+1}, Ae_j) - \sum_{k=1}^n \lambda_k (e_{n+1}, e_k) (e_k, e_j) \\ &= \lambda_j (e_{n+1}, e_j) - \lambda_j (e_{n+1}, e_j) = 0. \end{aligned}$$

To verify the second part of this claim,

$$\lambda_{n+1} e_{n+1} = A_n e_{n+1} = Ae_{n+1} - \sum_{k=1}^n \lambda_k e_k (e_{n+1}, e_k) = Ae_{n+1}$$

This proves the claim.

Claim 2: $|\lambda_n| \geq |\lambda_{n+1}|$.

Proof of claim: From 24.17 and the definition of A_n and $e_k \otimes e_k$,

$$(A_{n-1}e_{n+1}, e_{n+1}) = \left(\left(A - \sum_{k=1}^{n-1} \lambda_k e_k \otimes e_k \right) e_{n+1}, e_{n+1} \right) = (Ae_{n+1}, e_{n+1}) = (A_n e_{n+1}, e_{n+1})$$

Thus,

$$\lambda_{n+1} = (A_n e_{n+1}, e_{n+1}) = (A_{n-1} e_{n+1}, e_{n+1}) - \lambda_n |(e_n, e_{n+1})|^2 = (A_{n-1} e_{n+1}, e_{n+1})$$

By the previous claim. Therefore,

$$|\lambda_{n+1}| = |(A_{n-1} e_{n+1}, e_{n+1})| \leq |(A_{n-1} e_n, e_n)| = |\lambda_n|$$

by the definition of $|\lambda_n|$. (e_n makes $|(A_{n-1} x, x)|$ as large as possible.)

Claim 3: $\lim_{n \rightarrow \infty} \lambda_n = 0$.

Proof of claim: If for some n , $\lambda_n = 0$, then $\lambda_k = 0$ for all $k > n$ by claim 2. Thus, for some n , $A = \sum_{k=1}^n \lambda_k e_k \otimes e_k$. Assume then that $\lambda_k \neq 0$ for any k . Then if $\lim_{k \rightarrow \infty} |\lambda_k| = \varepsilon > 0$, one contradicts, $\|e_k\| = 1$ for all k because

$$\|Ae_n - Ae_m\|^2 = \|\lambda_n e_n - \lambda_m e_m\|^2 = \lambda_n^2 + \lambda_m^2 \geq 2\varepsilon^2$$

which shows there is no Cauchy subsequence of $\{Ae_n\}_{n=1}^\infty$, which contradicts the compactness of A . This proves the claim.

Claim 4: $\|A_n\| \rightarrow 0$

Proof of claim: Let $x, y \in B$

$$\begin{aligned} |\lambda_{n+1}| &\geq \left| \left(A_n \frac{x+y}{2}, \frac{x+y}{2} \right) \right| = \left| \frac{1}{4} (A_n x, x) + \frac{1}{4} (A_n y, y) + \frac{1}{2} (A_n x, y) \right| \\ &\geq \frac{1}{2} |(A_n x, y)| - \frac{1}{4} |(A_n x, x) + (A_n y, y)| \\ &\geq \frac{1}{2} |(A_n x, y)| - \frac{1}{4} (|(A_n x, x)| + |(A_n y, y)|) \geq \frac{1}{2} |(A_n x, y)| - \frac{1}{2} |\lambda_{n+1}| \end{aligned}$$

and so $3|\lambda_{n+1}| \geq |(A_n x, y)|$. It follows $\|A_n\| \leq 3|\lambda_{n+1}|$. By 24.17 this proves 24.16 and completes the proof. ■

24.3.2 Measurable Compact Operators

Here the operators will be of the form $A(\omega)$ where $\omega \in \Omega$ and $\omega \rightarrow A(\omega)x$ is strongly measurable and $A(\omega)$ is a compact operator in $\mathcal{L}(H, H)$.

Theorem 24.3.10 *Let $A(\omega) \in \mathcal{L}(H, H)$ be a compact self adjoint operator and H is a separable Hilbert space such that $\omega \rightarrow A(\omega)x$ is strongly measurable. Then there exist real numbers $\{\lambda_k(\omega)\}_{k=1}^\infty$ and vectors $\{e_k(\omega)\}_{k=1}^\infty$ such that*

$$\|e_k(\omega)\| = 1$$

$$\begin{aligned}
(e_k(\omega), e_j(\omega))_H &= 0 \text{ if } k \neq j, \\
A(\omega) e_k(\omega) &= \lambda_k(\omega) e_k(\omega), \\
|\lambda_n(\omega)| &\geq |\lambda_{n+1}(\omega)| \text{ for all } n, \\
\lim_{n \rightarrow \infty} \lambda_n(\omega) &= 0, \\
\lim_{n \rightarrow \infty} \left\| A(\omega) - \sum_{k=1}^n \lambda_k(\omega) (e_k(\omega) \otimes e_k(\omega)) \right\|_{\mathcal{L}(H, H)} &= 0.
\end{aligned}$$

The function $\omega \rightarrow \lambda_j(\omega)$ is measurable and $\omega \rightarrow e_j(\omega)$ is strongly measurable.

Proof: It is simply a repeat of the above proof of the Hilbert Schmidt theorem except at every step when the e_k and λ_k are defined, you use the Kuratowski measurable selection theorem, Theorem 24.3.4 on Page 663 to obtain $\lambda_k(\omega)$ is measurable and that $\omega \rightarrow e_k(\omega)$ is also measurable. This follows because the closed unit ball in a separable Hilbert space is a compact metric space.

When you consider $\max_{x \in B} |(A_n(\omega)x, x)|$, let $\psi(x, \omega) = |(A_n(\omega)x, x)|$. Then ψ is continuous in x by Lemma 24.3.6 on Page 664 and it is measurable in ω by assumption. Therefore, by the Kuratowski theorem, $e_k(\omega)$ is measurable in the sense that inverse images of weakly open sets in B are measurable. However, by Lemma 24.1.12 on Page 651 this is the same as weakly measurable. Since H is separable, this implies $\omega \rightarrow e_k(\omega)$ is also strongly measurable. The measurability of λ_k and e_k is the only new thing here and so this completes the proof. ■

24.4 Fubini's Theorem for Bochner Integrals

Now suppose $(\Omega_1, \mathcal{F}, \mu)$ and $(\Omega_2, \mathcal{S}, \lambda)$ are two σ finite measure spaces. Recall the notion of product measure. There was a σ algebra, denoted by $\mathcal{F} \times \mathcal{S}$ which is the smallest σ algebra containing the elementary sets, (finite disjoint unions of measurable rectangles) and a measure, denoted by $\mu \times \lambda$ defined on this σ algebra such that for $E \in \mathcal{F} \times \mathcal{S}$,

$$s_1 \rightarrow \lambda(E_{s_1}), (E_{s_1} \equiv \{s_2 : (s_1, s_2) \in E\})$$

is μ measurable and

$$s_2 \rightarrow \mu(E_{s_2}), (E_{s_2} \equiv \{s_1 : (s_1, s_2) \in E\})$$

is λ measurable. In terms of nonnegative functions which are $\mathcal{F} \times \mathcal{S}$ measurable,

$$\begin{aligned}
s_1 &\rightarrow f(s_1, s_2) \text{ is } \mu \text{ measurable,} \\
s_2 &\rightarrow f(s_1, s_2) \text{ is } \lambda \text{ measurable,} \\
s_1 &\rightarrow \int_{\Omega_2} f(s_1, s_2) d\lambda \text{ is } \mu \text{ measurable,} \\
s_2 &\rightarrow \int_{\Omega_1} f(s_1, s_2) d\mu \text{ is } \lambda \text{ measurable,}
\end{aligned}$$

and the conclusion of Fubini's theorem holds.

$$\begin{aligned}
\int_{\Omega_1 \times \Omega_2} f d(\mu \times \lambda) &= \int_{\Omega_1} \int_{\Omega_2} f(s_1, s_2) d\lambda d\mu \\
&= \int_{\Omega_2} \int_{\Omega_1} f(s_1, s_2) d\mu d\lambda.
\end{aligned}$$

The following theorem is the version of Fubini's theorem valid for Bochner integrable functions.

Theorem 24.4.1 *Let $f : \Omega_1 \times \Omega_2 \rightarrow X$ be strongly measurable with respect to $\mu \times \lambda$ and suppose*

$$\int_{\Omega_1 \times \Omega_2} \|f(s_1, s_2)\| d(\mu \times \lambda) < \infty. \quad (24.18)$$

Then there exist a set of μ measure zero, N and a set of λ measure zero, M such that the following formula holds with all integrals making sense.

$$\begin{aligned} \int_{\Omega_1 \times \Omega_2} f(s_1, s_2) d(\mu \times \lambda) &= \int_{\Omega_1} \int_{\Omega_2} f(s_1, s_2) \mathcal{X}_{N^c}(s_1) d\lambda d\mu \\ &= \int_{\Omega_2} \int_{\Omega_1} f(s_1, s_2) \mathcal{X}_{M^c}(s_2) d\mu d\lambda. \end{aligned}$$

Proof: First note that from 24.18 and the usual Fubini theorem for nonnegative valued functions,

$$\int_{\Omega_1 \times \Omega_2} \|f(s_1, s_2)\| d(\mu \times \lambda) = \int_{\Omega_1} \int_{\Omega_2} \|f(s_1, s_2)\| d\lambda d\mu$$

and so

$$\int_{\Omega_2} \|f(s_1, s_2)\| d\lambda < \infty \quad (24.19)$$

for μ a.e. s_1 . Say for all $s_1 \notin N$ where $\mu(N) = 0$.

Let $\phi \in X'$. Then $\phi \circ f$ is $\mathcal{F} \times \mathcal{S}$ measurable and

$$\begin{aligned} &\int_{\Omega_1 \times \Omega_2} |\phi \circ f(s_1, s_2)| d(\mu \times \lambda) \\ &\leq \int_{\Omega_1 \times \Omega_2} \|\phi\| \|f(s_1, s_2)\| d(\mu \times \lambda) < \infty \end{aligned}$$

and so from the usual Fubini theorem for complex valued functions,

$$\int_{\Omega_1 \times \Omega_2} \phi \circ f(s_1, s_2) d(\mu \times \lambda) = \int_{\Omega_1} \int_{\Omega_2} \phi \circ f(s_1, s_2) d\lambda d\mu. \quad (24.20)$$

Now also if you fix s_2 , it follows from the definition of strongly measurable and the properties of product measure mentioned above that $s_1 \rightarrow f(s_1, s_2)$ is strongly measurable. Also, by 24.19 $\int_{\Omega_2} \|f(s_1, s_2)\| d\lambda < \infty$ for $s_1 \notin N$. Therefore, by Theorem 24.2.4 $s_2 \rightarrow f(s_1, s_2) \mathcal{X}_{N^c}(s_1)$ is Bochner integrable. By 24.20 and 24.6

$$\begin{aligned} &\int_{\Omega_1 \times \Omega_2} \phi \circ f(s_1, s_2) d(\mu \times \lambda) \\ &= \int_{\Omega_1} \int_{\Omega_2} \phi \circ f(s_1, s_2) d\lambda d\mu \\ &= \int_{\Omega_1} \int_{\Omega_2} \phi(f(s_1, s_2) \mathcal{X}_{N^c}(s_1)) d\lambda d\mu \\ &= \int_{\Omega_1} \phi \left(\int_{\Omega_2} f(s_1, s_2) \mathcal{X}_{N^c}(s_1) d\lambda \right) d\mu. \end{aligned} \quad (24.21)$$

Each iterated integral makes sense and

$$\begin{aligned} s_1 &\rightarrow \int_{\Omega_2} \phi(f(s_1, s_2) \mathcal{X}_{N^c}(s_1)) d\lambda \\ &= \phi\left(\int_{\Omega_2} f(s_1, s_2) \mathcal{X}_{N^c}(s_1) d\lambda\right) \end{aligned} \quad (24.22)$$

is μ measurable because

$$\begin{aligned} (s_1, s_2) &\rightarrow \phi(f(s_1, s_2) \mathcal{X}_{N^c}(s_1)) \\ &= \phi(f(s_1, s_2)) \mathcal{X}_{N^c}(s_1) \end{aligned}$$

is product measurable. Now consider the function,

$$s_1 \rightarrow \int_{\Omega_2} f(s_1, s_2) \mathcal{X}_{N^c}(s_1) d\lambda. \quad (24.23)$$

I want to show this is also Bochner integrable with respect to μ so I can factor out ϕ once again. It's measurability follows from the Pettis theorem and the above observation 24.22. Also,

$$\begin{aligned} &\int_{\Omega_1} \left\| \int_{\Omega_2} f(s_1, s_2) \mathcal{X}_{N^c}(s_1) d\lambda \right\| d\mu \\ &\leq \int_{\Omega_1} \int_{\Omega_2} \|f(s_1, s_2)\| d\lambda d\mu \\ &= \int_{\Omega_1 \times \Omega_2} \|f(s_1, s_2)\| d(\mu \times \lambda) < \infty. \end{aligned}$$

Therefore, the function in 24.23 is indeed Bochner integrable and so in 24.21 the ϕ can be taken outside the last integral. Thus,

$$\begin{aligned} &\phi\left(\int_{\Omega_1 \times \Omega_2} f(s_1, s_2) d(\mu \times \lambda)\right) \\ &= \int_{\Omega_1 \times \Omega_2} \phi \circ f(s_1, s_2) d(\mu \times \lambda) \\ &= \int_{\Omega_1} \int_{\Omega_2} \phi \circ f(s_1, s_2) d\lambda d\mu \\ &= \int_{\Omega_1} \phi\left(\int_{\Omega_2} f(s_1, s_2) \mathcal{X}_{N^c}(s_1) d\lambda\right) d\mu \\ &= \phi\left(\int_{\Omega_1} \int_{\Omega_2} f(s_1, s_2) \mathcal{X}_{N^c}(s_1) d\lambda d\mu\right). \end{aligned}$$

Since X' separates the points,

$$\int_{\Omega_1 \times \Omega_2} f(s_1, s_2) d(\mu \times \lambda) = \int_{\Omega_1} \int_{\Omega_2} f(s_1, s_2) \mathcal{X}_{N^c}(s_1) d\lambda d\mu.$$

The other formula follows from similar reasoning. ■

24.5 The Spaces $L^p(\Omega; X)$

Recall that x is Bochner when it is strongly measurable and $\int_{\Omega} \|x(s)\| d\mu < \infty$. It is natural to generalize to $\int_{\Omega} \|x(s)\|^p d\mu < \infty$.

Definition 24.5.1 $x \in L^p(\Omega; X)$ for $p \in [1, \infty)$ if x is strongly measurable and

$$\int_{\Omega} \|x(s)\|^p d\mu < \infty$$

Also

$$\|x\|_{L^p(\Omega; X)} \equiv \|x\|_p \equiv \left(\int_{\Omega} \|x(s)\|^p d\mu \right)^{1/p}. \quad (24.24)$$

As in the case of scalar valued functions, two functions in $L^p(\Omega; X)$ are considered equal if they are equal a.e. With this convention, and using the same arguments found in the presentation of scalar valued functions it is clear that $L^p(\Omega; X)$ is a normed linear space with the norm given by 24.24. In fact, $L^p(\Omega; X)$ is a Banach space. This is the main contribution of the next theorem.

Lemma 24.5.2 If x_n is a Cauchy sequence in $L^p(\Omega; X)$ satisfying

$$\sum_{n=1}^{\infty} \|x_{n+1} - x_n\|_p < \infty,$$

then there exists $x \in L^p(\Omega; X)$ such that $x_n(s) \rightarrow x(s)$ a.e. and

$$\|x - x_n\|_p \rightarrow 0.$$

Proof: Let $g_N(s) \equiv \sum_{n=1}^N \|x_{n+1}(s) - x_n(s)\|_X$. Then by the triangle inequality,

$$\begin{aligned} \left(\int_{\Omega} g_N(s)^p d\mu \right)^{1/p} &\leq \sum_{n=1}^N \left(\int_{\Omega} \|x_{n+1}(s) - x_n(s)\|^p d\mu \right)^{1/p} \\ &\leq \sum_{n=1}^{\infty} \|x_{n+1} - x_n\|_p < \infty. \end{aligned}$$

Let

$$g(s) = \lim_{N \rightarrow \infty} g_N(s) = \sum_{n=1}^{\infty} \|x_{n+1}(s) - x_n(s)\|_X.$$

By the monotone convergence theorem,

$$\left(\int_{\Omega} g(s)^p d\mu \right)^{1/p} = \lim_{N \rightarrow \infty} \left(\int_{\Omega} g_N(s)^p d\mu \right)^{1/p} < \infty.$$

Therefore, there exists a measurable set of measure 0 called E , such that for $s \notin E$, $g(s) < \infty$. Hence, for $s \notin E$, $\lim_{N \rightarrow \infty} x_{N+1}(s)$ exists because

$$x_{N+1}(s) = x_{N+1}(s) - x_1(s) + x_1(s) = \sum_{n=1}^N (x_{n+1}(s) - x_n(s)) + x_1(s).$$

Thus, if $N > M$, and s is a point where $g(s) < \infty$,

$$\begin{aligned} \|x_{N+1}(s) - x_{M+1}(s)\|_X &\leq \sum_{n=M+1}^N \|x_{n+1}(s) - x_n(s)\|_X \\ &\leq \sum_{n=M+1}^{\infty} \|x_{n+1}(s) - x_n(s)\|_X \end{aligned}$$

which shows that $\{x_{N+1}(s)\}_{N=1}^{\infty}$ is a Cauchy sequence for each $s \notin E$. Now let

$$x(s) \equiv \begin{cases} \lim_{N \rightarrow \infty} x_N(s) & \text{if } s \notin E, \\ 0 & \text{if } s \in E. \end{cases}$$

Theorem 24.1.10 shows that x is strongly measurable. By Fatou's lemma,

$$\int_{\Omega} \|x(s) - x_N(s)\|^p d\mu \leq \liminf_{M \rightarrow \infty} \int_{\Omega} \|x_M(s) - x_N(s)\|^p d\mu.$$

But if N and M are large enough with $M > N$,

$$\left(\int_{\Omega} \|x_M(s) - x_N(s)\|^p d\mu \right)^{1/p} \leq \sum_{n=N}^M \|x_{n+1} - x_n\|_p \leq \sum_{n=N}^{\infty} \|x_{n+1} - x_n\|_p < \varepsilon$$

and this shows, since ε is arbitrary, that

$$\lim_{N \rightarrow \infty} \int_{\Omega} \|x(s) - x_N(s)\|^p d\mu = 0.$$

It remains to show $x \in L^p(\Omega; X)$. This follows from the above and the triangle inequality. Thus, for N large enough,

$$\begin{aligned} \left(\int_{\Omega} \|x(s)\|^p d\mu \right)^{1/p} &\leq \left(\int_{\Omega} \|x_N(s)\|^p d\mu \right)^{1/p} \\ &+ \left(\int_{\Omega} \|x(s) - x_N(s)\|^p d\mu \right)^{1/p} \leq \left(\int_{\Omega} \|x_N(s)\|^p d\mu \right)^{1/p} + \varepsilon < \infty. \blacksquare \end{aligned}$$

Theorem 24.5.3 $L^p(\Omega; X)$ is complete. Also every Cauchy sequence has a subsequence which converges pointwise.

Proof: If $\{x_n\}$ is Cauchy in $L^p(\Omega; X)$, extract a subsequence $\{x_{n_k}\}$ satisfying

$$\|x_{n_{k+1}} - x_{n_k}\|_p \leq 2^{-k}$$

and apply Lemma 24.5.2. The pointwise convergence of this subsequence was established in the proof of this lemma. This proves the theorem because if a subsequence of a Cauchy sequence converges, then the Cauchy sequence must also converge. \blacksquare

Observation 24.5.4 If the measure space is Lebesgue measure then you have continuity of translation in $L^p(\mathbb{R}^n; X)$ in the usual way. More generally, for μ a Radon measure on Ω a locally compact Hausdorff space, $C_c(\Omega; X)$ is dense in $L^p(\Omega; X)$. Here $C_c(\Omega; X)$ is the space of continuous X valued functions which have compact support in Ω . The proof of this little observation follows immediately from approximating with simple functions and then applying the appropriate considerations to the simple functions.

Clearly Fatou's lemma and the monotone convergence theorem make no sense for functions with values in a Banach space but the dominated convergence theorem holds in this setting.

Theorem 24.5.5 *If x is strongly measurable and $x_n(s) \rightarrow x(s)$ a.e. (for s off a set of measure zero) with*

$$\|x_n(s)\| \leq g(s) \text{ a.e.}$$

where $\int_{\Omega} g d\mu < \infty$, then x is Bochner integrable and

$$\int_{\Omega} x(s) d\mu = \lim_{n \rightarrow \infty} \int_{\Omega} x_n(s) d\mu.$$

Proof: The measurability of x follows from Theorem 24.1.10 if convergence happens for each s . Otherwise, x is measurable by assumption. Then $\|x_n(s) - x(s)\| \leq 2g(s)$ a.e. so, from Fatou's lemma,

$$\begin{aligned} \int_{\Omega} 2g(s) d\mu &\leq \liminf_{n \rightarrow \infty} \int_{\Omega} (2g(s) - \|x_n(s) - x(s)\|) d\mu \\ &= \int_{\Omega} 2g(s) d\mu - \limsup_{n \rightarrow \infty} \int_{\Omega} \|x_n(s) - x(s)\| d\mu \end{aligned}$$

and so,

$$\limsup_{n \rightarrow \infty} \int_{\Omega} \|x_n(s) - x(s)\| d\mu \leq 0$$

Also, from Fatou's lemma again,

$$\int_{\Omega} \|x(s)\| d\mu \leq \liminf_{n \rightarrow \infty} \int_{\Omega} \|x_n(s)\| d\mu < \int_{\Omega} g(s) d\mu < \infty$$

so $x \in L^1$. Then by the triangle inequality,

$$\limsup_{n \rightarrow \infty} \left\| \int_{\Omega} x(s) d\mu - \int_{\Omega} x_n(s) d\mu \right\| \leq \limsup_{n \rightarrow \infty} \int_{\Omega} \|x_n(s) - x(s)\| d\mu = 0 \blacksquare$$

One can also give a version of the Vitali convergence theorem.

Definition 24.5.6 *Let $\mathcal{A} \subseteq L^1(\Omega; X)$. Then \mathcal{A} is said to be uniformly integrable if for every $\varepsilon > 0$ there exists $\delta > 0$ such that whenever $\mu(E) < \delta$, it follows*

$$\int_E \|f\|_X d\mu < \varepsilon$$

for all $f \in \mathcal{A}$. It is bounded if

$$\sup_{f \in \mathcal{A}} \int_{\Omega} \|f\|_X d\mu < \infty.$$

Theorem 24.5.7 *Let $(\Omega, \mathcal{F}, \mu)$ be a finite measure space and let X be a separable Banach space. Let $\{f_n\} \subseteq L^1(\Omega; X)$ be uniformly integrable and bounded such that $f_n(\omega) \rightarrow f(\omega)$ for each $\omega \in \Omega$. Then $f \in L^1(\Omega; X)$ and*

$$\lim_{n \rightarrow \infty} \int_{\Omega} \|f_n - f\|_X d\mu = 0.$$

Proof: Let $\varepsilon > 0$ be given. Then by uniform integrability there exists $\delta > 0$ such that if $\mu(E) < \delta$ then

$$\int_E \|f_n\| d\mu < \varepsilon/3.$$

By Fatou's lemma the same inequality holds for f . Fatou's lemma shows $f \in L^1(\Omega; X)$, f being measurable because of Theorem 9.1.2.

By Egoroff's theorem, Theorem 24.1.18, there exists a set of measure less than δ , E such that the convergence of $\{f_n\}$ to f is uniform off E . Therefore,

$$\begin{aligned} \int_{\Omega} \|f - f_n\| d\mu &\leq \int_E (\|f\|_X + \|f_n\|_X) d\mu + \int_{E^c} \|f - f_n\|_X d\mu \\ &< \frac{2\varepsilon}{3} + \int_{E^c} \frac{\varepsilon}{(\mu(\Omega) + 1)3} d\mu < \varepsilon \end{aligned}$$

if n is large enough. ■

Note that a convenient way to achieve uniform integrability is to say $\{f_n\}$ is bounded in $L^p(\Omega; X)$ for some $p > 1$. This follows from Holder's inequality.

$$\int_E \|f_n\| d\mu \leq \left(\int_E d\mu \right)^{1/p'} \left(\int_{\Omega} \|f_n\|^p d\mu \right)^{1/p} \leq C \mu(E)^{1/p'}.$$

The following theorem is interesting.

Theorem 24.5.8 *Let $1 \leq p < \infty$ and let $p < r \leq \infty$. Then $L^r([0, T], X)$ is a Borel subset of $L^p([0, T], X)$. Letting $C([0, T], X)$ denote the continuous functions having values in X , $C([0, T], X)$ is also a Borel subset of $L^p([0, T], X)$. Here the measure is ordinary one dimensional Lebesgue measure on $[0, T]$.*

Proof: First consider the claim about $L^r([0, T], X)$. Let

$$B_M \equiv \left\{ x \in L^p([0, T], X) : \|x\|_{L^r([0, T], X)} \leq M \right\}.$$

Then B_M is a closed subset of $L^p([0, T], X)$. Here is why. If $\{x_n\}$ is a sequence of elements of B_M and $x_n \rightarrow x$ in $L^p([0, T], X)$, then passing to a subsequence, still denoted by x_n , it can be assumed $x_n(s) \rightarrow x(s)$ a.e. Hence Fatou's lemma can be applied to conclude

$$\int_0^T \|x(s)\|^r ds \leq \liminf_{n \rightarrow \infty} \int_0^T \|x_n(s)\|^r ds \leq M^r < \infty.$$

Now $\cup_{M=1}^{\infty} B_M = L^r([0, T], X)$. Note this did not depend on the measure space used. It would have been equally valid on any measure space.

Consider now $C([0, T], X)$. The norm on this space is the usual norm, $\|\cdot\|_{\infty}$. The argument above shows $\|\cdot\|_{\infty}$ is a Borel measurable function on $L^p([0, T], X)$. This is because $B_M \equiv \{x \in L^p([0, T], X) : \|x\|_{\infty} \leq M\}$ is a closed, hence Borel subset of $L^p([0, T], X)$. Now let $\theta \in \mathcal{L}(L^p([0, T], X), L^p(\mathbb{R}; X))$ such that $\theta(x(t)) = x(t)$ for all $t \in [0, T]$ and also $\theta \in \mathcal{L}(C([0, T], X), BC(\mathbb{R}; X))$ where $BC(\mathbb{R}; X)$ denotes the bounded continuous functions with a norm given by $\|x\| \equiv \sup_{t \in \mathbb{R}} \|x(t)\|$, and θx has compact support.

For example, you could define

$$\tilde{x}(t) \equiv \begin{cases} x(t) & \text{if } t \in [0, T] \\ x(2T - t) & \text{if } t \in [T, 2T] \\ x(-t) & \text{if } t \in [-T, 0] \\ 0 & \text{if } t \notin [-T, 2T] \end{cases}$$

and let $\Phi \in C_c^\infty(-T, 2T)$ such that $\Phi(t) = 1$ for $t \in [0, T]$. Then you could let

$$\theta x(t) \equiv \Phi(t) \tilde{x}(t).$$

Then let $\{\phi_n\}$ be a mollifier and define

$$\psi_n x(t) \equiv \phi_n * \theta x(t).$$

It follows $\psi_n x$ is uniformly continuous because

$$\begin{aligned} & \|\psi_n x(t) - \psi_n x(t')\|_X \\ & \leq \int_{\mathbb{R}} |\phi_n(t' - s) - \phi_n(t - s)| \|\theta x(s)\|_X ds \\ & \leq C \|x\|_p \left(\int_{\mathbb{R}} |\phi_n(t' - s) - \phi_n(t - s)|^{p'} ds \right)^{1/p'} \end{aligned}$$

Also for $x \in C([0, T]; X)$, it follows from usual mollifier arguments that

$$\|\psi_n x - x\|_{L^\infty([0, T]; X)} \rightarrow 0.$$

Here is why. For $t \in [0, T]$,

$$\begin{aligned} \|\psi_n x(t) - x(t)\|_X & \leq \int_{\mathbb{R}} \phi_n(s) \|\theta x(t - s) - \theta x(t)\| ds \\ & \leq C_\theta \int_{-1/n}^{1/n} \phi_n(s) ds \varepsilon = C_\theta \varepsilon \end{aligned}$$

provided n is large enough due to the compact support and consequent uniform continuity of θx .

If $\|\psi_n x - x\|_{L^\infty([0, T]; X)} \rightarrow 0$, then $\{\psi_n x\}$ is a Cauchy sequence in $C([0, T]; X)$ and this requires that x equals a continuous function a.e. Thus $C([0, T]; X)$ consists exactly of those functions, x of $L^p([0, T]; X)$ such that $\|\psi_n x - x\|_\infty \rightarrow 0$. It follows

$$\begin{aligned} C([0, T]; X) &= \\ & \cap_{n=1}^\infty \cup_{m=1}^\infty \cap_{k=m}^\infty \left\{ x \in L^p([0, T]; X) : \|\psi_k x - x\|_\infty \leq \frac{1}{n} \right\}. \end{aligned} \quad (24.25)$$

It only remains to show

$$S \equiv \{x \in L^p([0, T]; X) : \|\psi_k x - x\|_\infty \leq \alpha\}$$

is a Borel set. Suppose then that $x_n \in S$ and $x_n \rightarrow x$ in $L^p([0, T]; X)$. Then there exists a subsequence, still denoted by n such that $x_n \rightarrow x$ pointwise a.e. as well as in L^p . There exists a set of measure 0 such that for all n , and t not in this set,

$$\begin{aligned} \|\psi_k x_n(t) - x_n(t)\| & \equiv \left\| \int_{-1/k}^{1/k} \phi_k(s) (\theta x_n(t - s)) ds - x_n(t) \right\| \leq \alpha \\ x_n(t) & \rightarrow x(t). \end{aligned}$$

Then

$$\begin{aligned}
& \|\psi_k x_n(t) - x_n(t) - (\psi_k x(t) - x(t))\| \\
& \leq \|x_n(t) - x(t)\|_X + \left\| \int_{-1/k}^{1/k} \phi_k(s) (\theta x_n(t-s) - \theta x(t-s)) ds \right\| \\
& \leq \|x_n(t) - x(t)\|_X + C_{k,\theta} \|x_n - x\|_{L^p(0,T;X)}
\end{aligned}$$

which converges to 0 as $n \rightarrow \infty$. It follows that for a.e. t ,

$$\|\psi_k x(t) - x(t)\| \leq \alpha.$$

Thus S is closed and so the set in 24.25 is a Borel set. ■

As in the scalar case, the following lemma holds in this more general context.

Lemma 24.5.9 *Let (Ω, μ) be a regular measure space where Ω is a locally compact Hausdorff space or more simply a metric space with closed balls compact. Then $C_c(\Omega; X)$ the space of continuous functions having compact support and values in X is dense in $L^p(0, T; X)$ for all $p \in [0, \infty)$. For any measure space $(\Omega, \mathcal{F}, \mu)$, the simple functions are dense in $L^p(0, T; X)$.*

Proof: First, the simple functions are dense in $L^p(0, T; X)$. Let $f \in L^p(0, T; X)$ and let $\{x_n\}$ denote a sequence of simple functions which converge to f pointwise which also have the property that

$$\|x_n(s)\| \leq 2\|f(s)\|$$

Then

$$\int_{\Omega} \|x_n(s) - f(s)\|^p d\mu \rightarrow 0$$

from the dominated convergence theorem. Therefore, the simple functions are indeed dense in $L^p(0, T; X)$.

Next suppose (Ω, μ) is a regular measure space. If $x(s) \equiv \sum_i a_i \mathcal{X}_{E_i}(s)$ is a simple function, then by regularity, there exist compact sets K_i and open sets, V_i such that $K_i \subseteq E_i \subseteq V_i$ and $\mu(V_i \setminus K_i)^{1/p} < \varepsilon / \sum_i \|a_i\|$. Let $K_i \prec h_i \prec V_i$. Then consider

$$\sum_i a_i h_i \in C_c(\Omega).$$

By the triangle inequality,

$$\begin{aligned}
& \left(\int_{\Omega} \left\| \sum_i a_i h_i(s) - a_i \mathcal{X}_{E_i}(s) \right\|^p d\mu \right)^{1/p} \\
& \leq \sum_i \left(\int_{\Omega} \|a_i(h_i(s) - \mathcal{X}_{E_i}(s))\|^p d\mu \right)^{1/p} \\
& \leq \sum_i \left(\int_{\Omega} \|a_i\|^p |h_i(s) - \mathcal{X}_{E_i}(s)|^p d\mu \right)^{1/p} \leq \sum_i \|a_i\| \left(\int_{V_i \setminus K_i} d\mu \right)^{1/p} \\
& \leq \sum_i \|a_i\| \mu(V_i \setminus K_i)^{1/p} < \varepsilon
\end{aligned}$$

This and the first part of the lemma shows that $\overline{C_c(\Omega; X)} = L^p(\Omega; X)$. ■

24.6 Measurable Representatives

In this section consider the special case where $X = L^1(B, \nu)$ where (B, \mathcal{F}, ν) is a measure space and $x \in L^1(\Omega; X)$. Thus for each $s \in \Omega$, $x(s) \in L^1(B, \nu)$. In general, the map

$$(s, t) \rightarrow x(s)(t)$$

will not be product measurable, but one can obtain a measurable representative. This is important because it allows the use of Fubini's theorem on the measurable representative.

By Theorem 24.2.4, there exists a sequence of simple functions, $\{x_n\}$, of the form

$$x_n(s) = \sum_{k=1}^m a_k \mathcal{X}_{E_k}(s) \quad (24.26)$$

where $a_k \in L^1(B, \nu)$ which satisfy the conditions of Definition 24.2.3 and

$$\|x_n - x_m\|_{L^1(\Omega, L^1(B))} \rightarrow 0 \text{ as } m, n \rightarrow \infty \quad (24.27)$$

For such a simple function, you can assume the E_k are disjoint and then

$$\begin{aligned} \|x_n\|_{L^1(\Omega, L^1(B))} &= \sum_{k=1}^m \|a_k\|_{L^1(B)} \mu(E_k) = \sum_{k=1}^m \int_B |a_k| d\nu \mu(E_k) \\ &= \int_{\Omega} \int_B |a_k(t)| d\nu(t) \mathcal{X}_{E_k}(s) d\mu(s) \\ &= \int_{\Omega} \int_B |x_n| d\nu d\mu \end{aligned}$$

Also, each x_n is product measurable. Thus from 24.27,

$$\|x_n - x_m\|_{L^1(\Omega, L^1(B))} = \int_{\Omega} \int_B |x_n - x_m| d\nu d\mu$$

which shows that $\{x_n\}$ is a Cauchy sequence in $L^1(\Omega \times B, \mu \times \lambda)$. Then there exists $y \in L^1(\Omega \times B, \mu \times \lambda)$ and a subsequence still called $\{x_n\}$ such that

$$\begin{aligned} \lim_{n \rightarrow \infty} \int_{\Omega} \int_B |x_n - y| d\nu d\mu &= \lim_{n \rightarrow \infty} \int_{\Omega} \|x_n - y\|_{L^1(B)} d\mu \\ &= \|x_n - y\|_{L^1(\Omega, L^1(B))} = 0. \end{aligned}$$

Now consider 24.27. Since $\lim_{m \rightarrow \infty} x_m(s) = x(s)$ in $L^1(B)$, it follows from Fatou's lemma that

$$\|x_n - x\|_{L^1(\Omega, L^1(B))} \leq \liminf_{m \rightarrow \infty} \|x_n - x_m\|_{L^1(\Omega, L^1(B))} < \varepsilon$$

for all n large enough. Hence

$$\lim_{n \rightarrow \infty} \|x_n - x\|_{L^1(\Omega, L^1(B))} = 0$$

and so

$$x(s) = y(s) \text{ in } L^1(B) \text{ } \mu \text{ a.e. } s$$

In particular, for a.e. s , it follows that

$$x(s)(t) = y(s, t) \text{ for a.e. } t.$$

Now $\int_{\Omega} x(s) d\mu \in X = L^1(B, \nu)$ so it makes sense to ask for $(\int_{\Omega} x(s) d\mu)(t)$, at least μ a.e. t . To find what this is, note

$$\left\| \int_{\Omega} x_n(s) d\mu - \int_{\Omega} x(s) d\mu \right\|_X \leq \int_{\Omega} \|x_n(s) - x(s)\|_X d\mu.$$

Therefore, since the right side converges to 0,

$$\begin{aligned} \lim_{n \rightarrow \infty} \left\| \int_{\Omega} x_n(s) d\mu - \int_{\Omega} x(s) d\mu \right\|_X &= \\ \lim_{n \rightarrow \infty} \int_B \left| \left(\int_{\Omega} x_n(s) d\mu \right)(t) - \left(\int_{\Omega} x(s) d\mu \right)(t) \right| d\nu &= 0. \end{aligned}$$

But

$$\left(\int_{\Omega} x_n(s) d\mu \right)(t) = \int_{\Omega} x_n(s, t) d\mu \text{ a.e. } t.$$

Therefore

$$\lim_{n \rightarrow \infty} \int_B \left| \int_{\Omega} x_n(s, t) d\mu - \left(\int_{\Omega} x(s) d\mu \right)(t) \right| d\nu = 0. \quad (24.28)$$

Also, since $x_n \rightarrow y$ in $L^1(\Omega \times B)$,

$$\begin{aligned} 0 &= \lim_{n \rightarrow \infty} \int_B \int_{\Omega} |x_n(s, t) - y(s, t)| d\mu d\nu \geq \\ \lim_{n \rightarrow \infty} \int_B \left| \int_{\Omega} x_n(s, t) d\mu - \int_{\Omega} y(s, t) d\mu \right| d\nu &. \end{aligned} \quad (24.29)$$

From 24.28 and 24.29

$$\int_{\Omega} y(s, t) d\mu = \left(\int_{\Omega} x(s) d\mu \right)(t) \text{ a.e. } t.$$

Thus the following theorem is obtained.

Theorem 24.6.1 *Let $X = L^1(B)$ where (B, \mathcal{F}, ν) is a σ finite measure space and let $x \in L^1(\Omega; X)$. Then there exists a measurable representative, $y \in L^1(\Omega \times B)$, such that*

$$x(s) = y(s, \cdot) \text{ a.e. } s \text{ in } \Omega, \text{ the equation in } L^1(B),$$

and

$$\int_{\Omega} y(s, t) d\mu = \left(\int_{\Omega} x(s) d\mu \right)(t) \text{ a.e. } t.$$

24.7 Vector Measures

There is also a concept of vector measures.

Definition 24.7.1 *Let (Ω, \mathcal{S}) be a set and a σ algebra of subsets of Ω . A mapping*

$$F : \mathcal{S} \rightarrow X$$

is said to be a vector measure if

$$F(\cup_{i=1}^{\infty} E_i) = \sum_{i=1}^{\infty} F(E_i)$$

whenever $\{E_i\}_{i=1}^{\infty}$ is a sequence of disjoint elements of \mathcal{S} . For F a vector measure,

$$|F|(A) \equiv \sup \left\{ \sum_{F \in \pi(A)} \|F(F)\| : \pi(A) \text{ is a partition of } A \right\}.$$

This is the same definition that was given in the case where F would have values in \mathbb{C} , the only difference being the fact that now F has values in a general Banach space X as the vector space of values of the vector measure. Recall that a partition of A is a finite set, $\{F_1, \dots, F_m\} \subseteq \mathcal{S}$ such that $\cup_{i=1}^m F_i = A$. The same theorem about $|F|$ proved in the case of complex valued measures holds in this context with the same proof. For completeness, it is included here.

Theorem 24.7.2 *If $|F|(\Omega) < \infty$, then $|F|$ is a measure on \mathcal{S} .*

Proof: Let E_1 and E_2 be sets of \mathcal{S} such that $E_1 \cap E_2 = \emptyset$ and let $\{A_1^i, \dots, A_{n_i}^i\} = \pi(E_i)$, a partition of E_i which is chosen such that

$$|F|(E_i) - \varepsilon < \sum_{j=1}^{n_i} \|F(A_j^i)\| \quad i = 1, 2.$$

Consider the sets which are contained in either of $\pi(E_1)$ or $\pi(E_2)$, it follows this collection of sets is a partition of $E_1 \cup E_2$ which is denoted here by $\pi(E_1 \cup E_2)$. Then by the above inequality and the definition of total variation,

$$|F|(E_1 \cup E_2) \geq \sum_{F \in \pi(E_1 \cup E_2)} \|F(F)\| > |F|(E_1) + |F|(E_2) - 2\varepsilon,$$

which shows that since $\varepsilon > 0$ was arbitrary,

$$|F|(E_1 \cup E_2) \geq |F|(E_1) + |F|(E_2). \quad (24.30)$$

Let $\{E_j\}_{j=1}^{\infty}$ be a sequence of disjoint sets of \mathcal{S} and let $E_{\infty} = \cup_{j=1}^{\infty} E_j$. Then by the definition of total variation there exists a partition of E_{∞} , $\pi(E_{\infty}) = \{A_1, \dots, A_n\}$ such that

$$|F|(E_{\infty}) - \varepsilon < \sum_{i=1}^n \|F(A_i)\|.$$

Also,

$$A_i = \cup_{j=1}^{\infty} A_i \cap E_j, \text{ so } F(A_i) = \sum_{j=1}^{\infty} F(A_i \cap E_j)$$

and so by the triangle inequality, $\|F(A_i)\| \leq \sum_{j=1}^{\infty} \|F(A_i \cap E_j)\|$. Therefore, by the above,

$$|F|(E_{\infty}) - \varepsilon < \sum_{i=1}^n \overbrace{\sum_{j=1}^{\infty} \|F(A_i \cap E_j)\|}^{\geq \|F(A_i)\|} = \sum_{j=1}^{\infty} \sum_{i=1}^n \|F(A_i \cap E_j)\| \leq \sum_{j=1}^{\infty} |F|(E_j)$$

because $\{A_i \cap E_j\}_{i=1}^n$ is a partition of E_j .

Since $\varepsilon > 0$ is arbitrary, this shows

$$|F|(\cup_{j=1}^{\infty} E_j) \leq \sum_{j=1}^{\infty} |F|(E_j).$$

Also, 24.30 implies that whenever the E_i are disjoint, $|F|(\cup_{j=1}^n E_j) \geq \sum_{j=1}^n |F|(E_j)$. Therefore,

$$\sum_{j=1}^{\infty} |F|(E_j) \geq |F|(\cup_{j=1}^{\infty} E_j) \geq |F|(\cup_{j=1}^n E_j) \geq \sum_{j=1}^n |F|(E_j).$$

Since n is arbitrary,

$$|F|(\cup_{j=1}^{\infty} E_j) = \sum_{j=1}^{\infty} |F|(E_j)$$

which shows that $|F|$ is a measure as claimed. ■

Definition 24.7.3 A Banach space is said to have the Radon Nikodym property if whenever

$(\Omega, \mathcal{S}, \mu)$ is a finite measure space

$F : \mathcal{S} \rightarrow X$ is a vector measure with $|F|(\Omega) < \infty$

$$F \ll \mu$$

then one may conclude there exists $g \in L^1(\Omega; X)$ such that

$$F(E) = \int_E g(s) d\mu$$

for all $E \in \mathcal{S}$.

Some Banach spaces have the Radon Nikodym property and some don't. No attempt is made to give a complete answer to the question of which Banach spaces have this property, but the next theorem gives examples of many spaces which do. This next lemma was used earlier. I am presenting it again.

Lemma 24.7.4 Suppose ν is a complex measure defined on \mathcal{S} a σ algebra where (Ω, \mathcal{S}) is a measurable space, and let μ be a measure on \mathcal{S} with $|\nu(E)| \leq r\mu(E)$ and suppose there is $h \in L^1(\Omega, \mu)$ such that for all $E \in \mathcal{S}$,

$$\nu(E) = \int_E h d\mu,$$

Then $|h| \leq r$ a.e.

Proof: Let $B(p, \delta) \subseteq \mathbb{C} \setminus \overline{B(0, r)}$ and let $E \equiv h^{-1}(B(p, \delta))$. If $\mu(E) > 0$. Then

$$\left| \frac{1}{\mu(E)} \int_E h d\mu - p \right| \leq \frac{1}{\mu(E)} \int_E |h(\omega) - p| d\mu < \delta$$

Thus, $\left| \frac{\nu(E)}{\mu(E)} - p \right| < \delta$ and so $|\nu(E) - p\mu(E)| < \delta\mu(E)$ which implies

$$|\nu(E)| \geq (|p| - \delta)\mu(E) > r\mu(E) \geq |\nu(E)|$$

which contradicts the assumption. Hence $h^{-1}(B(p, \delta))$ is a set of μ measure zero for all such balls contained in $\mathbb{C} \setminus \overline{B(0, r)}$ and so, since countably many of these balls cover $\mathbb{C} \setminus \overline{B(0, r)}$, it follows that $\mu\left(h^{-1}\left(\mathbb{C} \setminus \overline{B(0, r)}\right)\right) = 0$ and so $|h(\omega)| \leq r$ for a.e. ω . ■

Theorem 24.7.5 *Suppose X' is a separable dual space. Then X' has the Radon Nikodym property.*

Proof: By Theorem 24.1.16, X is separable. Let D be a countable dense subset of X . Let $F \ll \mu$, μ a finite measure and F a vector measure and let $|F|(\Omega) < \infty$. Pick $x \in X$ and consider the map

$$E \rightarrow F(E)(x)$$

for $E \in \mathcal{S}$. This defines a complex measure which is absolutely continuous with respect to $|F|$. Therefore, by the earlier Radon Nikodym theorem, there exists $f_x \in L^1(\Omega, |F|)$ such that

$$F(E)(x) = \int_E f_x(s) d|F|. \quad (24.31)$$

Also, by definition $\|F(E)\| \leq |F|(E)$ so $|F(E)(x)| \leq |F|(E)\|x\|$. By Lemma 24.7.4, $|f_x(s)| \leq \|x\|$ for $|F|$ a.e. s . Let \tilde{D} consist of all finite linear combinations of the form $\sum_{i=1}^m a_i x_i$ where a_i is a rational point of \mathbb{F} and $x_i \in D$. For each of these countably many vectors, there is an exceptional set of measure zero off which $|f_x(s)| \leq \|x\|$. Let N be the union of all of them and define $f_x(s) \equiv 0$ if $s \notin N$. Then since $F(E)$ is in X' , it is linear and so for $\sum_{i=1}^m a_i x_i \in \tilde{D}$,

$$\begin{aligned} \int_E f_{\sum_{i=1}^m a_i x_i}(s) d|F| &= F(E) \left(\sum_{i=1}^m a_i x_i \right) = \sum_{i=1}^m a_i F(E)(x_i) \\ &= \int_E \sum_{i=1}^m a_i f_{x_i}(s) d|F| \end{aligned}$$

and so by uniqueness in the Radon Nikodym theorem,

$$f_{\sum_{i=1}^m a_i x_i}(s) = \sum_{i=1}^m a_i f_{x_i}(s) \quad |F| \text{ a.e.}$$

and so, we can regard this as holding for all $s \notin N$. Also, if $x \in \tilde{D}$, $|f_x(s)| \leq \|x\|$. Now for $x, y \in \tilde{D}$,

$$|f_x(s) - f_y(s)| = |f_{x-y}(s)| \leq \|x - y\|$$

and so, by density of \tilde{D} , we can define

$$h_x(s) \equiv \lim_{n \rightarrow \infty} f_{x_n}(s) \text{ where } x_n \rightarrow x, x_n \in \tilde{D}$$

For $s \in N$, all functions equal 0. Thus for all x , $|h_x(s)| \leq \|x\|$. The dominated convergence theorem and continuity of $F(E)$ implies that for $x_n \rightarrow x$, with $x_n \in \tilde{D}$,

$$\int_E h_x(s) d|F| = \lim_{n \rightarrow \infty} \int_E f_{x_n}(s) d|F| = \lim_{n \rightarrow \infty} F(E)(x_n) = F(E)(x). \quad (24.32)$$

It follows from the density of \tilde{D} that for all $x, y \in X$, $s \notin N$, and $a, b \in \mathbb{F}$, let $x_n \rightarrow x, y_n \rightarrow y, a_n \rightarrow a, b_n \rightarrow b$, with $x_n, y_n \in \tilde{D}$ and $a_n, b_n \in \mathbb{Q}$ or $\mathbb{Q} + i\mathbb{Q}$ in case $\mathbb{F} = \mathbb{C}$. Then

$$h_{ax+by}(s) = \lim_{n \rightarrow \infty} f_{a_n x_n + b_n y_n}(s) = \lim_{n \rightarrow \infty} a_n f_{x_n}(s) + b_n f_{y_n}(s) \equiv ah_x(s) + bh_y(s), \quad (24.33)$$

Let $\theta(s)$ be given by $\theta(s)(x) = h_x(s)$ if $s \notin N$ and let $\theta(s) = 0$ if $s \in N$. By 24.33 it follows that $\theta(s) \in X'$ for each s . Also

$$\theta(s)(x) = h_x(s) \in L^1(\Omega)$$

so $\theta(\cdot)$ is weak $*$ measurable. Since X' is separable, Theorem 24.1.15 implies that θ is strongly measurable. Furthermore, by 24.33,

$$\|\theta(s)\| \equiv \sup_{\|x\| \leq 1} |\theta(s)(x)| \leq \sup_{\|x\| \leq 1} |h_x(s)| \leq 1.$$

Therefore, $\int_{\Omega} \|\theta(s)\| d|F| < \infty$ so $\theta \in L^1(\Omega; X')$. Thus, if $E \in \mathcal{S}$,

$$\int_E h_x(s) d|F| = \int_E \theta(s)(x) d|F| = \left(\int_E \theta(s) d|F| \right)(x). \quad (24.34)$$

From 24.32 and 24.34, $(\int_E \theta(s) d|F|)(x) = F(E)(x)$ for all $x \in X$ and therefore,

$$\int_E \theta(s) d|F| = F(E).$$

Finally, since $F \ll \mu$, $|F| \ll \mu$ also and so there exists $k \in L^1(\Omega)$ such that

$$|F|(E) = \int_E k(s) d\mu$$

for all $E \in \mathcal{S}$, by the scalar Radon Nikodym Theorem. It follows

$$F(E) = \int_E \theta(s) d|F| = \int_E \theta(s) k(s) d\mu.$$

Letting $g(s) = \theta(s)k(s)$, this has proved the theorem. ■

Since each reflexive Banach spaces is a dual space, the following corollary holds.

Corollary 24.7.6 *Any separable reflexive Banach space has the Radon Nikodym property.*

It is not necessary to assume separability in the above corollary. For the proof of a more general result, consult *Vector Measures* by Diestel and Uhl, [13].

24.8 The Riesz Representation Theorem

The Riesz representation theorem for the spaces $L^p(\Omega; X)$ holds under certain conditions. The proof follows the proofs given earlier for scalar valued functions.

Definition 24.8.1 *If X and Y are two Banach spaces, X is isometric to Y if there exists $\theta \in \mathcal{L}(X, Y)$ such that*

$$\|\theta x\|_Y = \|x\|_X.$$

This will be written as $X \cong Y$. The map θ is called an isometry.

The next theorem says that $L^{p'}(\Omega; X')$ is always isometric to a subspace of $(L^p(\Omega; X))'$ for any Banach space, X .

Theorem 24.8.2 *Let X be any Banach space and let $(\Omega, \mathcal{S}, \mu)$ be a measure space. Then for $p > 1$, $L^{p'}(\Omega; X')$ is isometric to a subspace of $(L^p(\Omega; X))'$. Also, for $g \in L^{p'}(\Omega; X')$,*

$$\sup_{\|f\|_p \leq 1} \left| \int_{\Omega} g(\omega)(f(\omega)) d\mu \right| = \|g\|_{p'}.$$

If $p = 1$ and $p' = \infty$, this is still true assuming $\mu(\Omega) < \infty$.

Proof: First observe that for $f \in L^p(\Omega; X)$ and $g \in L^{p'}(\Omega; X')$,

$$\omega \rightarrow g(\omega)(f(\omega))$$

is a function in $L^1(\Omega)$. (To obtain measurability, write f as a limit of simple functions. Holder's inequality then yields the function is in $L^1(\Omega)$.) Define

$$\theta : L^{p'}(\Omega; X') \rightarrow (L^p(\Omega; X))'$$

by

$$\theta g(f) \equiv \int_{\Omega} g(\omega)(f(\omega)) d\mu.$$

Holder's inequality implies

$$\|\theta g\| \leq \|g\|_{p'} \quad (24.35)$$

and it is also clear that θ is linear. Next it is required to show $\|\theta g\| = \|g\|_{p'}$. To begin with, always assume $p > 1$.

This will first be verified for simple functions. Assume $\|g\|_{p'} \neq 0$ since if not, there is nothing to show. Let

$$g(\omega) = \sum_{i=1}^m c_i^* \mathcal{X}_{E_i}(\omega), \quad g \in L^{p'}(\Omega; X'), \quad c_i^* \neq 0$$

where $0 \neq c_i^* \in X'$, the E_i are disjoint. Let $d_i \in X$ be such that $\|d_i\|_X = 1$ and

$$c_i^*(d_i) \geq \|c_i^*\| - \varepsilon$$

Then let

$$f(\omega) \equiv \frac{1}{\|g\|_{L^{p'}(\Omega; X')}^{p'-1}} \sum_{i=1}^m d_i \|c_i^*\|^{p'-1} \mathcal{X}_{E_i}(\omega)$$

Then since $p' - 1 = p'/p$,

$$\int_{\Omega} \|f\|^p d\mu = \frac{1}{\|g\|_{L^{p'}(\Omega; X')}^{p'}} \int_{\Omega} \sum_{i=1}^m \|c_i^*\|^{p'} \mathcal{X}_{E_i}(\omega) d\mu = \frac{\|g\|_{L^{p'}(\Omega; X')}^{p'}}{\|g\|_{L^{p'}(\Omega; X')}^{p'}} = 1$$

Also $\int_{\Omega} g(\omega)(f(\omega)) =$

$$\begin{aligned} & \frac{1}{\|g\|_{L^{p'}(\Omega; X')}} \int_{\Omega} \sum_{i=1}^m c_i^*(d_i) \|c_i^*\|^{p'-1} \mathcal{X}_{E_i}(\omega) d\mu \\ & \geq \frac{1}{\|g\|_{L^{p'}(\Omega; X')}} \int_{\Omega} \sum_{i=1}^m (\|c_i^*\| - \varepsilon) \|c_i^*\|^{p'-1} \mathcal{X}_{E_i}(\omega) d\mu \\ & = \frac{1}{\|g\|_{L^{p'}(\Omega; X')}} \sum_{i=1}^m \|c_i^*\|^{p'} \mu(E_i) - \frac{1}{\|g\|_{L^{p'}(\Omega; X')}} \varepsilon \sum_{i=1}^m \|c_i^*\|^{p'-1} \mu(E_i) \\ & = \|g\|_{L^{p'}(\Omega; X')} - \varepsilon \end{aligned}$$

Therefore, $\|g\|_{L^{p'}(\Omega; X')} \geq \|\theta(g)\| \geq |\int_{\Omega} g(\omega)(f(\omega))| \geq \|g\|_{L^{p'}(\Omega; X')} - \varepsilon$ and since ε is arbitrary, it follows that $\|g\|_{L^{p'}(\Omega; X')} = \|\theta(g)\|$ whenever g is a simple function.

In general, let $g \in L^{p'}(\Omega; X')$ and let g_n be a sequence of simple functions converging to g in $L^{p'}(\Omega; X')$. Such a sequence exists by Lemma 24.1.2. Let $g_n(\omega) \rightarrow g(\omega)$, $\|g_n(\omega)\| \leq 2\|g(\omega)\|$. Then each g_n is in $L^{p'}(\Omega; X')$ and by the dominated convergence theorem they converge to g in $L^{p'}(\Omega; X')$. Then for $\|\cdot\|$ the norm in $(L^p(\Omega; X))'$,

$$\|\theta g\| = \lim_{n \rightarrow \infty} \|\theta g_n\| = \lim_{n \rightarrow \infty} \|g_n\| = \|g\|.$$

This proves the theorem in case $p = 1$ and shows θ is the desired isometry.

Next suppose $p = 1$ and $g \in L^{\infty}(\Omega; X')$. It is still the case that $\|\theta g\| \leq \|g\|_{L^{\infty}(\Omega; X')}$. As above, one must choose f appropriately. In this case, assume μ is a finite measure. Begin with g a simple function $g(\omega) = \sum_{i=1}^m c_i^* \mathcal{X}_{E_i}(\omega)$. Suppose $\|c_1^*\|$ is at least as large as all other $\|c_i^*\|$ modify if the largest of these occurs at $k \neq 1$. Thus $\|g\|_{\infty} = \|c_1^*\|_{X'}$. Now let $c_1^*(d_1) \geq \|c_1^*\| - \varepsilon \mu(E_i)$, $\|d_1\|_X = 1$, and let $f(\omega) \equiv \frac{d_1}{\mu(E_i)} \mathcal{X}_{E_i}(\omega)$. Then $\int_{\Omega} \|f\| d\mu = 1$. Also

$$g(\omega)(f(\omega)) \frac{c_1^*(d_1)}{\mu(E_i)} \mathcal{X}_{E_i}(\omega) \geq \frac{1}{\mu(E_i)} \mathcal{X}_{E_i}(\omega) (\|c_1^*\| - \varepsilon \mu(E_i))$$

and so

$$\begin{aligned} |\theta g(f)| &= \left| \int_{\Omega} g(\omega)(f(\omega)) d\mu \right| \geq \left| \int_{\Omega} \left(\frac{1}{\mu(E_i)} \mathcal{X}_{E_i}(\omega) (\|c_1^*\| - \varepsilon \mu(E_i)) \right) d\mu \right| \\ &\geq \|c_1^*\| - \varepsilon \mu(\Omega) = \|g\|_{\infty} - \varepsilon \mu(\Omega) \end{aligned}$$

Thus

$$\|g\|_{\infty} \geq \|\theta g\| \geq \|g\|_{\infty} - \varepsilon \mu(\Omega)$$

and so, since ε is arbitrary, it follows that $\|\theta g\| = \|g\|_{L^{\infty}(\Omega; X')}$. Extending from simple functions to functions in $L^{\infty}(\Omega; X')$ goes as before. Approximate with simple functions and pass to a limit. ■

Theorem 24.8.3 *If X is a Banach space and X' has the Radon Nikodym property, then if $(\Omega, \mathcal{S}, \mu)$ is a finite measure space, $(L^p(\Omega; X))' \cong L^{p'}(\Omega; X')$ and in fact the mapping θ of Theorem 24.8.2 is onto.*

Proof: Let $l \in (L^p(\Omega; X))'$ and define $F(E) \in X'$ by $F(E)(x) \equiv l(\mathcal{X}_E(\cdot)x)$.

Lemma 24.8.4 F defined above is a vector measure with values in X' and $|F|(\Omega) < \infty$.

Proof of the lemma: Clearly $F(E)$ is linear. Also

$$\|F(E)\| = \sup_{\|x\| \leq 1} \|F(E)(x)\| \leq \|l\| \sup_{\|x\| \leq 1} \|\mathcal{X}_E(\cdot)x\|_{L^p(\Omega; X)} \leq \|l\| \mu(E)^{1/p}.$$

Let $\{E_i\}_{i=1}^\infty$ be a sequence of disjoint elements of \mathcal{S} and let $E = \cup_{n < \infty} E_n$.

$$\begin{aligned} \left| F(E)(x) - \sum_{k=1}^n F(E_k)(x) \right| &= \left| l(\mathcal{X}_E(\cdot)x) - \sum_{i=1}^n l(\mathcal{X}_{E_i}(\cdot)x) \right| \\ &\leq \|l\| \left\| \mathcal{X}_E(\cdot)x - \sum_{i=1}^n \mathcal{X}_{E_i}(\cdot)x \right\|_{L^p(\Omega; X)} \leq \|l\| \mu\left(\bigcup_{k>n} E_k\right)^{1/p} \|x\|. \end{aligned} \quad (24.36)$$

Since $\mu(\Omega) < \infty$, $\lim_{n \rightarrow \infty} \mu\left(\bigcup_{k>n} E_k\right)^{1/p} = 0$ and so inequality 24.36 shows that

$$\lim_{n \rightarrow \infty} \left\| F(E) - \sum_{k=1}^n F(E_k) \right\|_{X'} = 0.$$

To show $|F|(\Omega) < \infty$, let $\varepsilon > 0$ be given, let $\{H_1, \dots, H_n\}$ be a partition of Ω , and let $\|x_i\| \leq 1$ be chosen in such a way that $F(H_i)(x_i) > \|F(H_i)\| - \varepsilon/n$. Thus

$$\begin{aligned} -\varepsilon + \sum_{i=1}^n \|F(H_i)\| &< \sum_{i=1}^n |l(\mathcal{X}_{H_i}(\cdot)x_i)| \leq \|l\| \left\| \sum_{i=1}^n \mathcal{X}_{H_i}(\cdot)x_i \right\|_{L^p(\Omega; X)} \\ &\leq \|l\| \left(\int_{\Omega} \sum_{i=1}^n \mathcal{X}_{H_i}(s) d\mu \right)^{1/p} = \|l\| \mu(\Omega)^{1/p}. \end{aligned}$$

Since $\varepsilon > 0$ was arbitrary, $\sum_{i=1}^n \|F(H_i)\| < \|l\| \mu(\Omega)^{1/p}$. Since the partition was arbitrary, this shows $|F|(\Omega) \leq \|l\| \mu(\Omega)^{1/p}$ and this proves the lemma. ■

Continuing with the proof of Theorem 24.8.3, note that $F \ll \mu$. Since X' has the Radon Nikodym property, there exists $g \in L^1(\Omega; X')$ such that $F(E) = \int_E g(s) d\mu$. Also, from the definition of $F(E)$,

$$\begin{aligned} l\left(\sum_{i=1}^n x_i \mathcal{X}_{E_i}(\cdot)\right) &= \sum_{i=1}^n l(\mathcal{X}_{E_i}(\cdot)x_i) \\ &= \sum_{i=1}^n F(E_i)(x_i) = \sum_{i=1}^n \int_{E_i} g(s)(x_i) d\mu. \end{aligned} \quad (24.37)$$

It follows from 24.37 that whenever h is a simple function,

$$l(h) = \int_{\Omega} g(s)(h(s)) d\mu. \quad (24.38)$$

Let $G_n \equiv \{s : \|g(s)\|_{X'} \leq n\}$ and let $j : L^p(G_n; X) \rightarrow L^p(\Omega; X)$ be given by

$$jh(s) = \begin{cases} h(s) & \text{if } s \in G_n, \\ 0 & \text{if } s \notin G_n. \end{cases}$$

Letting h be a simple function in $L^p(G_n; X)$,

$$j^*l(h) = l(jh) = \int_{G_n} g(s)(h(s)) d\mu. \quad (24.39)$$

Since the simple functions are dense in $L^p(G_n; X)$, and $g \in L^{p'}(G_n; X')$, it follows 24.39 holds for all $h \in L^p(G_n; X)$. By Theorem 24.8.2,

$$\|g\|_{L^{p'}(G_n; X')} = \|j^*l\|_{(L^p(G_n; X))'} \leq \|l\|_{(L^p(\Omega; X))'}.$$

By the monotone convergence theorem, $\|g\|_{L^{p'}(\Omega; X')} = \lim_{n \rightarrow \infty} \|g\|_{L^{p'}(G_n; X')} \leq \|l\|_{(L^p(\Omega; X))'}$. Therefore $g \in L^{p'}(\Omega; X')$ and since simple functions are dense in $L^p(\Omega; X)$, 24.38 holds for all $h \in L^p(\Omega; X)$. Thus $l = \theta g$ and the theorem is proved because, by Theorem 24.8.2, $\|l\| = \|g\|$ and so the mapping θ is onto because l was arbitrary. ■

As in the scalar case, everything generalizes to the case of σ finite measure spaces. The proof is almost identical.

Lemma 24.8.5 *Let $(\Omega, \mathcal{S}, \mu)$ be a σ finite measure space and let X be a Banach space such that X' has the Radon Nikodym property. Then there exists a measurable function, r such that $r(x) > 0$ for all x , such that $|r(x)| < M$ for all x , and $\int r d\mu < \infty$. For*

$$\Lambda \in (L^p(\Omega; X))', \quad p \geq 1,$$

there exists a unique $h \in L^{p'}(\Omega; X')$, $L^\infty(\Omega; X')$ if $p = 1$ such that $\Lambda f = \int h(f) d\mu$. Also $\|h\| = \|\Lambda\|$. ($\|h\| = \|h\|_{p'}$ if $p > 1$, $\|h\|_\infty$ if $p = 1$). Here $\frac{1}{p} + \frac{1}{p'} = 1$.

Proof: First suppose r exists as described. Also, to save on notation and to emphasize the similarity with the scalar case, denote the norm in the various spaces by $|\cdot|$. Define a new measure $\tilde{\mu}$, according to the rule

$$\tilde{\mu}(E) \equiv \int_E r d\mu. \quad (24.40)$$

Thus $\tilde{\mu}$ is a finite measure on \mathcal{S} . Now define a mapping, $\eta : L^p(\Omega; X, \mu) \rightarrow L^p(\Omega; X, \tilde{\mu})$ by $\eta f = r^{-\frac{1}{p}} f$. Then

$$\|\eta f\|_{L^p(\tilde{\mu})}^p = \int \left| r^{-\frac{1}{p}} f \right|^p r d\mu = \|f\|_{L^p(\mu)}^p$$

and so η is one to one and in fact preserves norms. I claim that also η is onto. To see this, let $g \in L^p(\Omega; X, \tilde{\mu})$ and consider the function, $r^{\frac{1}{p}} g$. Then

$$\int \left| r^{\frac{1}{p}} g \right|^p d\mu = \int |g|^p r d\mu = \int |g|^p d\tilde{\mu} < \infty$$

Thus $r^{\frac{1}{p}} g \in L^p(\Omega; X, \mu)$ and $\eta(r^{\frac{1}{p}} g) = g$ showing that η is onto as claimed. Thus η is one to one, onto, and preserves norms. Consider the diagram below which is descriptive of the situation in which η^* must be one to one and onto.

$$\begin{array}{ccccc} h, L^{p'}(\tilde{\mu}) & L^p(\tilde{\mu})', \tilde{\Lambda} & \xrightarrow{\eta^*} & L^p(\mu)', \Lambda \\ & & \eta & & \\ & L^p(\tilde{\mu}) & \xleftarrow{} & L^p(\mu) \end{array}$$

Then for $\Lambda \in L^p(\mu)'$, there exists a unique $\tilde{\Lambda} \in L^p(\tilde{\mu})'$ such that $\eta^* \tilde{\Lambda} = \Lambda$, $\|\tilde{\Lambda}\| = \|\Lambda\|$. By the Riesz representation theorem for finite measure spaces, there exists a unique $h \in L^{p'}(\tilde{\mu}) \equiv L^{p'}(\Omega; X', \tilde{\mu})$ which represents $\tilde{\Lambda}$ in the manner described in the Riesz representation theorem. Thus $\|h\|_{L^{p'}(\tilde{\mu})} = \|\tilde{\Lambda}\| = \|\Lambda\|$ and for all $f \in L^p(\mu)$,

$$\begin{aligned} \Lambda(f) &= \eta^* \tilde{\Lambda}(f) \equiv \tilde{\Lambda}(\eta f) = \int h(\eta f) d\tilde{\mu} = \int rh \left(r^{-\frac{1}{p}} f\right) d\mu \\ &= \int r^{\frac{1}{p'}} h f d\mu. \end{aligned}$$

Now

$$\int \left| r^{\frac{1}{p'}} h \right|^{p'} d\mu = \int |h|^{p'} r d\mu = \|h\|_{L^{p'}(\tilde{\mu})}^{p'} < \infty.$$

Thus $\left\| r^{\frac{1}{p'}} h \right\|_{L^{p'}(\mu)} = \|h\|_{L^{p'}(\tilde{\mu})} = \|\tilde{\Lambda}\| = \|\Lambda\|$ and represents Λ in the appropriate way. If $p = 1$, then $1/p' \equiv 0$. Now consider the existence of r . Since the measure space is σ finite, there exist $\{\Omega_n\}$ disjoint, each having positive measure and their union equals Ω . Then define

$$r(\omega) \equiv \sum_{n=1}^{\infty} \frac{1}{n^2} \mu(\Omega_n)^{-1} \chi_{\Omega_n}(\omega)$$

This proves the Lemma.

Theorem 24.8.6 (Riesz representation theorem) *Let $(\Omega, \mathcal{S}, \mu)$ be σ finite and let X' have the Radon Nikodym property. Then for $\Lambda \in (L^p(\Omega; X, \mu))'$, $p \geq 1$ there exists a unique $h \in L^q(\Omega, X', \mu)$, $L^\infty(\Omega, X', \mu)$ if $p = 1$ such that $\Lambda f = \int h(f) d\mu$. Also $\|h\| = \|\Lambda\|$. ($\|h\| = \|h\|_q$ if $p > 1$, $\|h\|_\infty$ if $p = 1$). Here $\frac{1}{p} + \frac{1}{q} = 1$.*

Proof: The above lemma gives the existence part of the conclusion of the theorem. Uniqueness is done as before.

Corollary 24.8.7 *If X' is separable, then for $(\Omega, \mathcal{S}, \mu)$ a σ finite measure space,*

$$(L^p(\Omega; X))' \cong L^{p'}(\Omega; X').$$

Corollary 24.8.8 *If X is separable and reflexive, then for $(\Omega, \mathcal{S}, \mu)$ a σ finite measure space,*

$$(L^p(\Omega; X))' \cong L^{p'}(\Omega; X').$$

Corollary 24.8.9 *If X is separable and reflexive and $(\Omega, \mathcal{S}, \mu)$ a σ finite measure space, then if $p \in (1, \infty)$, then $L^p(\Omega; X)$ is reflexive.*

Proof: This is just like the scalar valued case.

24.9 An Example of Polish Space

Here is an interesting example. Obviously $L^\infty(0, T, H)$ is not separable with the normed topology. However, bounded sets turn out to be metric spaces which are complete and separable. This is the next lemma. Recall that a Polish space is a complete separable metric space. In this example, H is a separable real Hilbert space or more generally a separable real Banach space.

Lemma 24.9.1 Let $B = \overline{B(0, L)}$ be a closed ball in $L^\infty(0, T, H)$. Then B is a Polish space with respect to the weak $*$ topology. The closure is taken with respect to the usual topology.

Proof: Let $\{z_k\}_{k=1}^\infty = X$ be a dense countable subspace in $L^1(0, T, H)$. You start with a dense countable set and then consider all finite linear combinations having coefficients in \mathbb{Q} . Then the metric on B is

$$d(f, g) \equiv \sum_{k=1}^{\infty} 2^{-k} \frac{|\langle f - g, z_k \rangle_{L^\infty, L^1}|}{1 + |\langle f - g, z_k \rangle_{L^\infty, L^1}|}$$

Is B complete? Suppose you have a Cauchy sequence $\{f_n\}$. This happens if and only if $\{\langle f_n, z_k \rangle\}_{n=1}^\infty$ is a Cauchy sequence for each k . Therefore, there exists

$$\xi(z_k) = \lim_{n \rightarrow \infty} \langle f_n, z_k \rangle.$$

Then for $a, b \in \mathbb{Q}$, and $z, w \in X$

$$\xi(az + bw) = \lim_{n \rightarrow \infty} \langle f_n, az + bw \rangle = \lim_{n \rightarrow \infty} a \langle f_n, z \rangle + b \langle f_n, w \rangle = a\xi(z) + b\xi(w)$$

showing that ξ is linear on X a dense subspace of $L^1(0, T, H)$. Is ξ bounded on this dense subspace with bound L ? For $z \in X$,

$$|\xi(z)| \equiv \lim_{n \rightarrow \infty} |\langle f_n, z \rangle| \leq \limsup_{n \rightarrow \infty} \|f_n\|_{L^\infty} \|z\|_{L^1} \leq L \|z\|_{L^1}$$

Hence ξ is also bounded on this dense subset of $L^1(0, T, H)$. Therefore, there is a unique bounded linear extension of ξ to all of $L^1(0, T, H)$ still denoted as ξ such that its norm in $L^1(0, T, H)'$ is no larger than L . It follows from the Riesz representation theorem that there exists a unique $f \in L^\infty(0, T, H)$ such that for all $w \in L^1(0, T, H)$, $\xi(w) = \langle f, w \rangle$ and $\|f\| \leq L$. This f is the limit of the Cauchy sequence $\{f_n\}$ in B . Thus B is complete.

Is B separable? Let $f \in B$. Let $\varepsilon > 0$ be given. Choose M such that $\sum_{k=M+1}^\infty 2^{-k} < \frac{\varepsilon}{4}$. Then the finite set $\{z_1, \dots, z_M\}$ is uniformly integrable. There exists $\delta > 0$ such that if $m(S) < \delta$, then $\int_S |z_k|_H dm < \left(\frac{\varepsilon}{4(1+\|f\|_{L^\infty})}\right)$. Then there is a sequence of simple functions $\{s_n\}$ which converge uniformly to f off a set of measure zero, N , $\|s_n\|_{L^\infty} \leq \|f\|_{L^\infty}$. By regularity of the measure, there exists a continuous function with compact support h_n such that $s_n = h_n$ off a set of measure no more than $\delta/4^n$ and also $\|h_n\|_{L^\infty} \leq \|f\|_{L^\infty}$. Then off a set of measure no more than $\frac{1}{3}\delta$, $h_n(r) \rightarrow f(r)$. Now by Eggorov's theorem and outer regularity, one can enlarge this exceptional set to obtain an open set S of measure no more than $\delta/2$ such that the convergence is uniform off this exceptional set. Thus f equals the uniform limit of continuous functions on S^C . Define

$$h(r) \equiv \begin{cases} \lim_{n \rightarrow \infty} h_n(r) = f(r) & \text{on } S^C \\ 0 & \text{on } S \setminus N \\ 0 & \text{on } N \end{cases}$$

Then $\|h\|_{L^\infty} \leq \|f\|_{L^\infty}$. Now consider $\bar{h} * \psi_m(r)$ where ψ_r is approximate identity.

$$\begin{aligned} \psi_m(t) &= \frac{1}{2} m \mathcal{X}_{[-1/m, 1/m]}(t), \quad \bar{h} * \psi_m(t) \\ &= \frac{1}{2} m \int_{-1/m}^{1/m} \bar{h}(t-s) ds = \frac{1}{2} m \int_{t-1/m}^{t+1/m} \bar{h}(s) ds \end{aligned}$$

where we define \bar{h} to be the 0 extension of \bar{h} off $[0, T]$. This is a continuous function of t . Also *a.e.t* is a Lebesgue point and so for *a.e.t*,

$$\left| \frac{1}{2} m \int_{t-1/m}^{t+1/m} \bar{h}(s) ds - \bar{h}(t) \right| \rightarrow 0$$

$$|\bar{h} * \psi_m(r)| \equiv \left| \int_{\mathbb{R}} \bar{h}(r-s) \psi_m(s) ds \right| \leq \|\bar{h}\|_{L^\infty} \leq \|f\|_{L^\infty}$$

Thus this continuous function is in $L^\infty(0, T, H)$. Letting $z = z_k \in L^1(0, T, H)$ be one of those defined above,

$$\left| \int_0^T \langle \bar{h} * \psi_m(t) - f(t), z(t) \rangle dt \right| \leq \int_0^T |\langle \bar{h} * \psi_m(t) - h(t), z(t) \rangle| dt$$

$$+ \int_0^T |\langle h(t) - f(t), z(t) \rangle| dt \quad (24.41)$$

for a.e. t , $\bar{h} * \psi_m(t) - h(t) \rightarrow 0$ and the integrand in the first integral in the above is bounded by $2\|f\|_{L^\infty}|z(t)|_H$ so by the dominated convergence theorem, as $m \rightarrow \infty$, the first integral converges to 0. As to the second, it is dominated by

$$\int_S |\langle h(t) - f(t), z(t) \rangle| dt \leq 2\|f\|_{L^\infty} \int_S |z(t)| dt < \frac{2\|f\|_{L^\infty} \varepsilon}{4(1 + \|f\|_{L^\infty})} \leq \frac{\varepsilon}{2}$$

Therefore, choosing m large enough so that the first integral on the right in 24.41 is less than $\frac{\varepsilon}{4}$ for each z_k for $k \leq M$, then for each of these,

$$\begin{aligned} d(f, \bar{h} * \psi_m) &\leq \frac{\varepsilon}{4} + \sum_{k=1}^M 2^{-k} \frac{(\varepsilon/4) + (\varepsilon/2)}{1 + ((\varepsilon/4) + (\varepsilon/2))} = \frac{\varepsilon}{4} + \sum_{k=1}^M 2^{-k} \frac{3}{4} \frac{\varepsilon}{\varepsilon + 1} \\ &\leq \frac{\varepsilon}{4} + \frac{3\varepsilon}{4} \sum_{k=1}^M 2^{-k} < \frac{\varepsilon}{4} + \frac{3\varepsilon}{4} = \varepsilon \end{aligned}$$

which appears to show that $C([0, T], H)$ is weak * dense in $L^\infty(0, T, H)$. However, this last space is obviously separable in terms of the norm topology. Let D be a countable dense subset of $C([0, T], H)$. For $f \in L^\infty(0, T, H)$ let $g \in C([0, T], H)$ such that $d(f, g) < \frac{\varepsilon}{4}$. Then let $h \in D$ be so close to g in $C([0, T], H)$ that

$$\sum_{k=1}^M 2^{-k} \frac{|\langle h - g, z_k \rangle_{L^\infty, L^1}|}{1 + |\langle h - g, z_k \rangle_{L^\infty, L^1}|} < \frac{\varepsilon}{2}$$

Then $d(f, h) \leq d(f, g) + d(g, h) < \frac{\varepsilon}{4} + \frac{\varepsilon}{2} + \frac{\varepsilon}{4} = \varepsilon$ It appears that D is dense in B in the weak * topology. ■

24.10 Weakly Convergent Sequences

There is an interesting little result which relates to weak limits in $L^2(\Gamma, E)$ for E a Banach space. I am not sure where to put this thing but think that this would be a good place for it. It obviously generalizes to L^p spaces.

Proposition 24.10.1 *Let E be a Banach space and let $\{u_n\}$ be a sequence in $L^2(\Gamma, E)$ and let $G(x)$ be a weakly compact set in E , and $u_n(x) \in G(x)$ a.e. for each n . Let $\limsup \{u_n(x)\}$ denote the set of all weak limits of subsequences of $\{u_n(x)\}$ and let $H(x)$ be the closure of the convex hull of $\limsup \{u_n(x)\}$. Then if $u_n \rightarrow u$ weakly in $L^2(\Gamma, E)$, then $u(x) \in H(x)$ for a.e. x .*

Proof: Let $H = \{w \in L^2(\Gamma, E) : w(x) \in H(x) \text{ a.e.}\}$. Then H is convex. If you have $w_i \in H$, then since each $H(x)$ is convex, it follows that $\lambda w_1(x) + (1 - \lambda)w_2(x) \in H$ for a.e. x and $\lambda \in [0, 1]$. Is H closed? Suppose you have $w_n \in H$ and $w_n \rightarrow w$ in $L^2(\Gamma, E)$. Then there is a subsequence such that pointwise convergence happens a.e. and so since H is closed, you have $w(x) \in H$ for a.e. x . Hence H is also weakly closed in $L^2(\Gamma, E)$. Thus if u is the weak limit of $\{u_n\}$ in $L^2(\Gamma, E)$, it must be the case that $u(x) \in H(x)$ a.e. ■

As a case of this which might be pretty interesting, suppose $G(x)$ is not just weakly compact but also convex. Then $H(x) = G(x)$ and you can say that $u(x) \in H(x)$ a.e. whenever it is a weak limit in $L^2(\Gamma, E)$ of functions u_n for which $u_n(x) \in G(x)$.

24.11 Some Embedding Theorems

The next lemma is a very useful little result which involves embeddings of Banach spaces.

Lemma 24.11.1 *Suppose $V \subseteq W$ and the injection map is compact, hence continuous. Suppose also that $W \subseteq U$ with continuous injection. Then for any $\varepsilon > 0$ there exists C_ε such that for all $v \in V$, $\|v\|_W \leq \varepsilon \|v\|_V + C_\varepsilon \|v\|_U$.*

Proof: Suppose not. Then there exists $\varepsilon > 0$ for which things don't work out. Thus there exists $v_n \in V$ such that $\|v_n\|_W > \varepsilon \|v_n\|_V + n \|v_n\|_U$. Dividing by $\|v_n\|_V$, it can also be assumed that $\|v_n\|_V = 1$. Thus $\|v_n\|_W > \varepsilon + n \|v_n\|_U$, and so $\|v_n\|_U \rightarrow 0$. However, v_n is contained in the closed unit ball of V which is, by assumption precompact in W . Hence, there exists a subsequence, still denoted as $\{v_n\}$ such that $v_n \rightarrow v$ in W . But it was just determined that $v = 0$ and so $0 \geq \limsup_{n \rightarrow \infty} (\varepsilon + n \|v_n\|_U) \geq \varepsilon$ which is a contradiction. ■

Recall the following definition, this time for the space of continuous functions defined on a compact set with values in a Banach space.

Definition 24.11.2 *Let $\mathcal{A} \subseteq C(K; V)$ where the last symbol denotes the continuous functions defined on a compact set $K \subseteq X$ a metric space having values in V a Banach space. Then \mathcal{A} is equicontinuous if for every $\varepsilon > 0$, there exists $\delta > 0$ such that for every $f \in \mathcal{A}$, if $d(x, y) < \delta$, then $\|f(x) - f(y)\|_V < \varepsilon$. Also $\mathcal{A} \subseteq C(K; V)$ is uniformly bounded means*

$$\sup_{f \in \mathcal{A}} \|f\|_{\infty, V} < \infty \text{ where } \|f\|_{\infty, V} \equiv \max_{x \in K} \|f(x)\|_V.$$

Here is a general version of the Ascoli Arzela theorem valid for Banach spaces.

Theorem 24.11.3 *Let $V \subseteq W \subseteq U$ where the injection map of V into W is compact and W embeds continuously into U , these being Banach spaces. Assume:*

1. $\mathcal{A} \subseteq C(K; U)$ where K is compact and \mathcal{A} is equicontinuous.
2. $\sup_{f \in \mathcal{A}} \|f\|_{\infty, V} < \infty$ where $\|f\|_{\infty, V} \equiv \max_{x \in K} \|f(x)\|_V$.

Then

1. $\mathcal{A} \subseteq C(K; W)$ and \mathcal{A} is equicontinuous into W
2. \mathcal{A} is pre-compact in $C(K; W)$. This means that $\overline{\mathcal{A}}$ is compact in $C(K; W)$.

Proof: Let $C \equiv \sup_{f \in \mathcal{A}} \|f\|_{\infty, V} < \infty$. Let $\varepsilon > 0$ be given. Then from Lemma 24.11.1,

$$\|f(x) - f(y)\|_W \leq \frac{\varepsilon}{5C} \|f(x) - f(y)\|_V + C_\varepsilon \|f(x) - f(y)\|_U \leq \frac{2\varepsilon}{5} + C_\varepsilon \|f(x) - f(y)\|_U$$

By equicontinuity in $C(K, U)$, there exists a $\delta > 0$ such that if $d(x, y) < \delta$, then for all $f \in \mathcal{A}$, $C_\varepsilon \|f(x) - f(y)\|_U < \frac{2\varepsilon}{5}$. Thus if $d(x, y) < \delta$, then $\|f(x) - f(y)\|_W < \varepsilon$ for all $f \in \mathcal{A}$.

It remains to verify that \mathcal{A} is pre-compact in $C(K; W)$. Since this space of continuous functions is complete, it suffices to verify that for all $\varepsilon > 0$, \mathcal{A} has an ε net. Suppose then that for some $\varepsilon > 0$ there is no ε net. Thus there is an infinite sequence $\{f_n\}$ for which $\|f_n - f_m\|_{\infty, W} \geq \varepsilon$ whenever $m \neq n$. There exists $\delta > 0$ such that if $d(x, y) < \delta$, then for all f_n , $\|f_n(x) - f_n(y)\|_W < \frac{\varepsilon}{5}$. Let $\{x_k\}_{k=1}^p$ be a $\delta/2$ net for K . This is where we use K is compact. By compactness of the embedding of V into W , there exists a further subsequence, still called $\{f_n\}$ such that each $\{f_n(x_k)\}_{n=1}^\infty$ converges, this for each x_k in that $\delta/2$ net. Thus there is a single N such that if $n > N$, then for all $m, n > N$, and $k \leq p$, $\|f_n(x_k) - f_m(x_k)\|_W < \frac{\varepsilon}{5}$. Now letting $x \in K$ be arbitrary, it is in $B(x_k, \delta/2)$ for some x_k . Therefore, for n, m larger than N ,

$$\begin{aligned} \|f_n(x) - f_m(x)\|_W &\leq \|f_n(x) - f_n(x_k)\|_W + \|f_n(x_k) - f_m(x_k)\|_W + \|f_m(x_k) - f_m(x)\|_W \\ &< \frac{\varepsilon}{5} + \frac{\varepsilon}{5} + \frac{\varepsilon}{5} = \frac{3\varepsilon}{5} \end{aligned}$$

Taking the maximum for all x , for $m, n > N$, $\|f_n - f_m\|_{W, \infty} \leq \frac{3\varepsilon}{5} < \varepsilon$ contrary to the assumption that every pair is further apart than ε . Thus \mathcal{A} is totally bounded so its closure would also be totally bounded and complete. In other words, \mathcal{A} is pre-compact in $C(K; W)$. ■

In the following theorem about compact subsets of an L^p space, the measure will be Lebesgue measure. It depends on the above version of the Ascoli Arzela theorem. First note the following which I will use when convenient. For $a, b \geq 0$, and $p \geq 1$, then by convexity of $\phi(t) = t^p$ for $t \geq 0$, $(a+b)^p \leq 2^{p-1}(a^p + b^p)$. Also, for such p , $(a+b)^{1/p} \leq a^{1/p} + b^{1/p}$. Usually the thing of interest in this theorem is the case where $V = W = U = \mathbb{R}$. However, the more general version to be presented is interesting I think. Of course closed and bounded sets are compact in \mathbb{R} so the usual case works as a special case of what is about to be presented.

Theorem 24.11.4 *Let $V \subseteq W \subseteq U$ where these are Banach spaces such that the injection map of V into W is compact and the injection map of W into U is continuous. Let Ω be an open set in \mathbb{R}^m and let \mathcal{A} be a bounded subset of $L^p(\Omega; V)$ and suppose that for all $\varepsilon > 0$, there exist a $\delta > 0$ such that if $|\mathbf{h}| < \delta$, then for \tilde{u} denoting the zero extension of u off Ω ,*

$$\int_{\mathbb{R}^m} \|\tilde{u}(\mathbf{x} + \mathbf{h}) - \tilde{u}(\mathbf{x})\|_U^p dx < \varepsilon^p \quad (24.42)$$

Suppose also that for each $\varepsilon > 0$ there exists an open set, $G_\varepsilon \subseteq \Omega$ such that $\overline{G_\varepsilon} \subseteq \Omega$ is compact and for all $u \in \mathcal{A}$,

$$\int_{\Omega \setminus \overline{G_\varepsilon}} \|u(\mathbf{x})\|_W^p dx < \varepsilon^p \quad (24.43)$$

Then \mathcal{A} is precompact in $L^p(\mathbb{R}^n; W)$.

Proof: Let $\infty > M \geq \sup_{u \in L^p(\Omega; V)} \|u\|_{L^p(\Omega; V)}^p$. Let $\{\psi_n\}$ be a mollifier with support in $B(0, 1/n)$. I need to show that \mathcal{A} has an η net in $L^p(\Omega; W)$ for every $\eta > 0$. Suppose for some $\eta > 0$ it fails to have an η net. Without loss of generality, let $\eta < 1$. Then by 24.43, it follows that for small enough $\varepsilon > 0$, $\mathcal{A}_\varepsilon \equiv \{u \mathcal{X}_{\overline{G_\varepsilon}} : u \in \mathcal{A}\}$ fails to have an $\eta/2$ net. Indeed, pick ε small enough that for all $u \in \mathcal{A}$, $\|u \mathcal{X}_{\overline{G_\varepsilon}} - u\|_{L^p(\Omega; W)} < \frac{\eta}{5}$. Then if $\left\{u_k \mathcal{X}_{\overline{G_\varepsilon}}\right\}_{k=1}^r$ is an $\eta/2$ net for \mathcal{A}_ε , so that $\cup_{k=1}^r B\left(u_k \mathcal{X}_{\overline{G_\varepsilon}}, \frac{\eta}{2}\right) \supseteq \mathcal{A}_\varepsilon$, then for $w \in \mathcal{A}$, $w \mathcal{X}_{\overline{G_\varepsilon}} \in B\left(u_k \mathcal{X}_{\overline{G_\varepsilon}}, \frac{\eta}{2}\right)$ for some u_k . Hence,

$$\begin{aligned} \|w - u_k\|_{L^p(\Omega; W)} &\leq \|w - w \mathcal{X}_{\overline{G_\varepsilon}}\|_{L^p(\Omega; W)} + \|w \mathcal{X}_{\overline{G_\varepsilon}} - u_k \mathcal{X}_{\overline{G_\varepsilon}}\|_{L^p(\Omega; W)} \\ &\quad + \|u_k \mathcal{X}_{\overline{G_\varepsilon}} - u_k\|_{L^p(\Omega; W)} \leq \frac{\eta}{5} + \frac{\eta}{2} + \frac{\eta}{5} < \eta \end{aligned}$$

and so $\{u_k\}_{k=1}^r$ would be an η net for \mathcal{A} which is assumed to not exist.

Pick this ε in all that follows. By compactness, Lemma 24.11.1, there exists C_η such that for all $u \in V$,

$$\|u\|_W^p \leq \frac{\eta}{50(2^{p-1})M} \|u\|_V^p + C_\eta \|u\|_U^p \quad (24.44)$$

Let $\mathcal{A}_{\varepsilon n}$ consist of $\mathcal{A}_{\varepsilon n} \equiv \{u \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n : u \in \mathcal{A}\}$. I want to show that $\mathcal{A}_{\varepsilon n}$ satisfies the conditions for Theorem 24.11.3.

Lemma 24.11.5 *For each n , $\mathcal{A}_{\varepsilon n}$ satisfies the conditions of Theorem 24.11.3.*

Proof: First consider the equicontinuity condition of that theorem. It suffices to show that if $\eta > 0$ then there exists $\delta > 0$ such that if $|\mathbf{h}| < \delta$, then for any $u \in \mathcal{A}$ and $\mathbf{x} \in \overline{G_\varepsilon}$,

$$\|u \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n(\mathbf{x} + \mathbf{h}) - u \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n(\mathbf{x})\|_U < \eta$$

Always assume $|\mathbf{h}| < \text{dist}(\overline{G_\varepsilon}, \Omega^C)$, and $\mathbf{x} \in \overline{G_\varepsilon}$. Also assume that $|\mathbf{h}|$ is small enough that

$$\begin{aligned} &\left(\int_{\mathbb{R}^m} \left| \left(\mathcal{X}_{\overline{G_\varepsilon}}(\mathbf{x} - \mathbf{y} + \mathbf{h}) - \mathcal{X}_{\overline{G_\varepsilon}}(\mathbf{x} - \mathbf{y}) \right) \psi_n(\mathbf{y}) \right|^{p'} dz \right)^{1/p'} = \\ &\left(\int_{\mathbb{R}^m} \left| \left(\mathcal{X}_{\overline{G_\varepsilon}}(\mathbf{z} + \mathbf{h}) - \mathcal{X}_{\overline{G_\varepsilon}}(\mathbf{z}) \right) \psi_n(\mathbf{x} - \mathbf{z}) \right|^{p'} dz \right)^{1/p'} < \frac{\eta}{2M} \end{aligned} \quad (24.45)$$

This can be obtained because by Holder's inequality,

$$\begin{aligned} &\left(\int_{\mathbb{R}^m} \left| \left(\mathcal{X}_{\overline{G_\varepsilon}}(\mathbf{z} + \mathbf{h}) - \mathcal{X}_{\overline{G_\varepsilon}}(\mathbf{z}) \right) \psi_n(\mathbf{x} - \mathbf{z}) \right|^{p'} dz \right)^{1/p'} \\ &\leq \left(\int_{\mathbb{R}^m} \left| \mathcal{X}_{\overline{G_\varepsilon}}(\mathbf{z} + \mathbf{h}) - \mathcal{X}_{\overline{G_\varepsilon}}(\mathbf{z}) \right|^{2p'} dz \right)^{\frac{1}{2p'}} \left(\int_{\mathbb{R}^m} \psi_n(\mathbf{x} - \mathbf{z})^{2p'} dz \right)^{\frac{1}{2p'}} \end{aligned}$$

which is small independent of \mathbf{x} for $|\mathbf{h}|$ small enough, thanks to continuity of translation in $L^{2p'}(\mathbb{R}^m)$. Then

$$\begin{aligned}
& \left\| u \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n(x + \mathbf{h}) - u \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n(x) \right\|_U \\
&= \left\| \int_{\mathbb{R}^m} \left(\tilde{u}(x + \mathbf{h} - \mathbf{y}) \mathcal{X}_{\overline{G_\varepsilon}}(x + \mathbf{h} - \mathbf{y}) - \tilde{u}(x - \mathbf{y}) \mathcal{X}_{\overline{G_\varepsilon}}(x - \mathbf{y}) \right) \psi_n(\mathbf{y}) d\mathbf{y} \right\|_U \\
&\leq \int_{\mathbb{R}^m} \left\| \left(\tilde{u}(x + \mathbf{h} - \mathbf{y}) \mathcal{X}_{\overline{G_\varepsilon}}(x + \mathbf{h} - \mathbf{y}) - \tilde{u}(x - \mathbf{y}) \mathcal{X}_{\overline{G_\varepsilon}}(x - \mathbf{y}) \right) \right\|_U \psi_n(\mathbf{y}) d\mathbf{y}
\end{aligned}$$

Changing the variables,

$$\begin{aligned}
&\leq \int_{\mathbb{R}^m} \left\| \begin{aligned} &(\tilde{u}(z + \mathbf{h}) - \tilde{u}(z)) \mathcal{X}_{\overline{G_\varepsilon}}(z + \mathbf{h}) \\ &+ \tilde{u}(z) \left(\mathcal{X}_{\overline{G_\varepsilon}}(z + \mathbf{h}) - \mathcal{X}_{\overline{G_\varepsilon}}(z) \right) \end{aligned} \right\|_U \psi_n(x - z) dz \\
&\leq \int_{\mathbb{R}^m} \left\| (\tilde{u}(z + \mathbf{h}) - \tilde{u}(z)) \mathcal{X}_{\overline{G_\varepsilon}}(z + \mathbf{h}) \right\|_U \psi_n(x - z) dz \\
&\quad + \int_{\mathbb{R}^m} \|\tilde{u}(z)\|_U \left| \mathcal{X}_{\overline{G_\varepsilon}}(z + \mathbf{h}) - \mathcal{X}_{\overline{G_\varepsilon}}(z) \right| \psi_n(x - z) dz \quad (24.46)
\end{aligned}$$

The first integral

$$\leq \left(\int_{\mathbb{R}^m} \|\tilde{u}(z + \mathbf{h}) - \tilde{u}(z)\|_U^p dz \right)^{1/p} \left(\int_{\mathbb{R}^m} \psi_n^{p'}(x - z) dz \right)^{1/p'}$$

You make the obvious change here in case $p = 1$. Instead of the above, you would have

$$\leq \int_{\mathbb{R}^m} \|\tilde{u}(z + \mathbf{h}) - \tilde{u}(z)\|_U dz 2 \|\psi_n\|_\infty$$

Since Lebesgue measure is translation independent, there is a constant C_n such that the above is $\leq C_n \left(\int_{\mathbb{R}^m} \|\tilde{u}(z + \mathbf{h}) - \tilde{u}(z)\|_U^p dz \right)^{1/p} < \eta/2$ and this holds for all $u \in \mathcal{A}$. As for the second integral in 24.46, from 24.45, it follows that this term is no larger than

$$\leq \left(\int_{\mathbb{R}^m} \|\tilde{u}(z)\|_U^p dz \right)^{1/p} \left(\int_{\mathbb{R}^m} \left(\left| \mathcal{X}_{\overline{G_\varepsilon}}(z + \mathbf{h}) - \mathcal{X}_{\overline{G_\varepsilon}}(z) \right| \psi_n(x - z) \right)^{p'} dz \right)^{1/p'}$$

and by 24.45, $< M \frac{\eta}{2M} = \frac{\eta}{2}$. Thus, if $\delta < \text{dist}(\overline{G_\varepsilon}, \Omega^C)$ and 24.45 holds, then for all $u \in \mathcal{A}$, when $|\mathbf{h}| < \delta$,

$$\left\| u \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n(x + \mathbf{h}) - u \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n(x) \right\|_U < \eta$$

and so the desired equicontinuity condition holds for $\mathcal{A}_{\varepsilon n}$. Note that δ does depend on n but for each n , things work out well.

I also need to verify that the functions in $\mathcal{A}_{\varepsilon n}$ are uniformly bounded. For $x \in \overline{G_\varepsilon}$ and $u \in \mathcal{A}$,

$$\begin{aligned}
\left\| u \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n(x) \right\|_V &\leq \int_{\overline{G_\varepsilon}} \|u(z)\| \psi_n(x - z) dz \\
&\leq \left(\int_{\Omega} \|u(z)\|^p dz \right)^{1/p} \left(\int_{\Omega} \psi_n(x - z)^{p'} dz \right)^{1/p'} \leq MC_n \blacksquare
\end{aligned}$$

Now is a general statement about norms, indicating that the L^p norm is no more than a constant times the norm involving the maximum.

$$\left(\int_{\overline{G_\varepsilon}} \|v(x)\|_W^p dx \right)^{1/p} \leq \max_{x \in \overline{G_\varepsilon}} \|v(x)\|_W m(\overline{G_\varepsilon}) \equiv m(\overline{G_\varepsilon}) \|v\|_{W, \infty}$$

It follows from Theorem 24.11.3 that for every $\eta > 0$, there exists a η net in $C(\overline{G_\varepsilon}; W)$ for $\mathcal{A}_{\varepsilon n}$, this for each n . Then from the above inequality, it follows that for each η , there exists an η net in $L^p(\overline{G_\varepsilon}; W)$ for $\mathcal{A}_{\varepsilon n}$.

Recall also, from the assumption that the theorem is not true, $\mathcal{A}_\varepsilon \equiv \{u \mathcal{X}_{\overline{G_\varepsilon}} : u \in \mathcal{A}\}$ has no $\eta/2$ net in $L^p(\overline{G_\varepsilon}; W)$. Next I estimate the distance in $L^p(\overline{G_\varepsilon}; W)$ between $u \mathcal{X}_{\overline{G_\varepsilon}}$ for $u \in \mathcal{A}$ and $u \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n$. The idea is that for each n , $\mathcal{A}_{\varepsilon n}$ has an $\eta/8$ net and for n large enough, $u \mathcal{X}_{\overline{G_\varepsilon}}$ is close to $u \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n$ so a contradiction will result if the functions of the second sort are totally bounded while those functions of the first sort don't. Assume always that $1/n < \text{dist}(\overline{G_\varepsilon}, \Omega^C)$. Using Minkowski's inequality,

$$\begin{aligned} & \left\| u \mathcal{X}_{\overline{G_\varepsilon}} - u \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n \right\|_{L^p(\overline{G_\varepsilon}; W)} = \\ & \left(\int_{\mathbb{R}^m} \left\| \int_{\mathbb{R}^m} (u \mathcal{X}_{\overline{G_\varepsilon}}(x) - u \mathcal{X}_{\overline{G_\varepsilon}}(x - y)) \psi_n(y) dy \right\|_W^p dx \right)^{1/p} \\ & \leq \int_{B(0, 1/n)} \psi_n(y) \left(\int_{\mathbb{R}^m} \left\| (u \mathcal{X}_{\overline{G_\varepsilon}}(x) - u \mathcal{X}_{\overline{G_\varepsilon}}(x - y)) \right\|_W^p dx \right)^{1/p} dy \\ & \leq \int_{B(0, 1/n)} \psi_n(y) \left(\int_{\mathbb{R}^m} \| \tilde{u}(x) - \tilde{u}(x - y) \|_W^p dx \right)^{1/p} dy \\ & \leq \int_{B(0, 1/n)} \psi_n(y) \left(\frac{\int_{\mathbb{R}^m} \frac{\eta}{50(2^{p-1})M} \| \tilde{u}(x) - \tilde{u}(x - y) \|_V^p dx}{+ C_\eta \int_{\mathbb{R}^m} \| \tilde{u}(x) - \tilde{u}(x - y) \|_U^p dx} \right)^{1/p} dy \\ & \leq \int_{B(0, \frac{1}{n})} \psi_n(y) \left(\frac{\int_{\mathbb{R}^m} \frac{\eta}{50(2^{p-1})M} 2^{p-1} 2 (\| \tilde{u}(x) \|_V^p) dx}{+ C_\eta \int_{\mathbb{R}^m} \| \tilde{u}(x) - \tilde{u}(x - y) \|_U^p dx} \right)^{1/p} dy \\ & \leq \int_{B(0, \frac{1}{n})} \psi_n(y) \left(\frac{\int_{\mathbb{R}^m} \frac{\eta}{25M} (\| \tilde{u}(x) \|_V^p) dx}{+ \int_{\mathbb{R}^m} C_\eta \| \tilde{u}(x) - \tilde{u}(x - y) \|_U^p dx} \right)^{1/p} dy \\ & \leq \int_{B(0, \frac{1}{n})} \psi_n(y) \left(\frac{\eta}{25} + \int_{\mathbb{R}^m} C_\eta \| \tilde{u}(x) - \tilde{u}(x - y) \|_U^p dx \right)^{1/p} dy \end{aligned}$$

By assumption 24.42, there exists N such that if $n \geq N$, then $|y| < \frac{1}{n}$ and for all $u \in \mathcal{A}$,

$$\begin{aligned} \left\| u \mathcal{X}_{\overline{G_\varepsilon}} - u \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n \right\|_{L^p(\overline{G_\varepsilon}; W)} & \leq \int_{B(0, \frac{1}{n})} \psi_n(y) \left(\frac{\eta}{25} + \frac{\eta^p}{8^p} \right)^{1/p} dy \\ & \leq \int_{B(0, \frac{1}{n})} \psi_n(y) \left(\frac{\eta}{25} + \frac{\eta}{8} \right) dy = \frac{\eta}{25} + \frac{\eta}{8} \end{aligned}$$

Recall $\eta < 1$.

Let n be this large. Then let $\left\{u_k \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n\right\}_{k=1}^r$ be a $\frac{\eta}{8}$ net for $\mathcal{A}_{\varepsilon n}$ in $L^p(\overline{G_\varepsilon}; W)$. Then consider the balls $B\left(u_k \mathcal{X}_{\overline{G_\varepsilon}}, \frac{\eta}{4}\right)$ in $L^p(\overline{G_\varepsilon}; W)$. If $w \mathcal{X}_{\overline{G_\varepsilon}}$ is in \mathcal{A}_ε , is it in some $B\left(u_k \mathcal{X}_{\overline{G_\varepsilon}}, \frac{\eta}{2}\right)$? By what was just shown, there is k such that

$$\left\|w \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n - u_k \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n\right\|_{L^p(\overline{G_\varepsilon}; W)} < \frac{\eta}{8}$$

and also

$$\left\|w \mathcal{X}_{\overline{G_\varepsilon}} - w \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n\right\|_{L^p(\overline{G_\varepsilon}; W)} < \frac{\eta}{8} + \frac{\eta}{25}$$

$$\left\|u_k \mathcal{X}_{\overline{G_\varepsilon}} - u_k \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n\right\|_{L^p(\overline{G_\varepsilon}; W)} < \frac{\eta}{8} + \frac{\eta}{25}$$

Thus,

$$\begin{aligned} & \left\|w \mathcal{X}_{\overline{G_\varepsilon}} - u_k \mathcal{X}_{\overline{G_\varepsilon}}\right\|_{L^p(\overline{G_\varepsilon}; W)} \leq \left\|w \mathcal{X}_{\overline{G_\varepsilon}} - w \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n\right\|_{L^p(\overline{G_\varepsilon}; W)} \\ & + \left\|w \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n - u_k \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n\right\|_{L^p(\overline{G_\varepsilon}; W)} + \left\|u_k \mathcal{X}_{\overline{G_\varepsilon}} * \psi_n - u_k \mathcal{X}_{\overline{G_\varepsilon}}\right\|_{L^p(\overline{G_\varepsilon}; W)} \\ & < \frac{3\eta}{8} + \frac{2\eta}{25} < \frac{\eta}{2} \end{aligned}$$

It follows that $\left\{u_k \mathcal{X}_{\overline{G_\varepsilon}}\right\}_{k=1}^r$ is a $\eta/2$ net for $L^p(\overline{G_\varepsilon}; W)$ contrary to the construction. Thus \mathcal{A} has an η net after all. ■

In case Ω is a closed interval, there are several versions of these sorts of embeddings which are enormously useful in the study of nonlinear evolution equations or inclusions.

The following theorem is an infinite dimensional version of the Ascoli Arzela theorem. It is like a well known result due to Simon [52]. It is an appropriate generalization when you do not necessarily have weak derivatives. I am giving another proof although Theorem 24.11.3 given above is actually more general.

Theorem 24.11.6 *Let $q > 1$ and let $E \subseteq W \subseteq X$ where the injection map is continuous from W to X and compact from E to W . Let S be defined by*

$$\left\{u \text{ such that } \|u(t)\|_E \leq R \text{ for all } t \in [a, b], \text{ and } \|u(s) - u(t)\|_X \leq R|t - s|^{1/q}\right\}.$$

Thus S is bounded in $L^\infty(a, b; E)$ and in addition, the functions are uniformly Hölder continuous into X . Then $S \subseteq C([a, b]; W)$ and if $\{u_n\} \subseteq S$, there exists a subsequence, $\{u_{n_k}\}$ which converges to a function $u \in C([a, b]; W)$ in the following way: $\lim_{k \rightarrow \infty} \|u_{n_k} - u\|_{\infty, W} = 0$.

Proof: First consider the issue of S being a subset of $C([a, b]; W)$. Let $\varepsilon > 0$ be given. Then by Lemma 24.11.1, there exists a constant, C_ε such that for all $u \in W$

$$\|u\|_W \leq \frac{\varepsilon}{6R} \|u\|_E + C_\varepsilon \|u\|_X.$$

Therefore, for all $u \in S$,

$$\|u(t) - u(s)\|_W \leq \frac{\varepsilon}{6R} \|u(t) - u(s)\|_E + C_\varepsilon \|u(t) - u(s)\|_X$$

$$\leq \frac{\varepsilon}{6R} (\|u(t)\|_E + \|u(s)\|_E) + C_\varepsilon \|u(t) - u(s)\|_X \leq \frac{\varepsilon}{3} + C_\varepsilon R |t - s|^{1/q}. \quad (24.47)$$

Since ε is arbitrary, it follows $u \in C([a, b]; W)$.

Let $D = \mathbb{Q} \cap [a, b]$ so D is a countable dense subset of $[a, b]$. Let $D = \{t_n\}_{n=1}^\infty$. By compactness of the embedding of E into W , there exists a subsequence $u_{(n,1)}$ such that as $n \rightarrow \infty$, $u_{(n,1)}(t_1)$ converges to a point in W . Now take a subsequence of this, called $(n, 2)$ such that as $n \rightarrow \infty$, $u_{(n,2)}(t_2)$ converges to a point in W . It follows that $u_{(n,2)}(t_1)$ also converges to a point of W . Continue this way. Now consider the diagonal sequence, $u_k \equiv u_{(k,k)}$. This sequence is a subsequence of $u_{(n,l)}$ whenever $k > l$. Therefore, $u_k(t_j)$ converges for all $t_j \in D$.

Claim: Let $\{u_k\}$ be as just defined, converging at every point of $D \equiv [a, b] \cap \mathbb{Q}$. Then $\{u_k\}$ converges at every point of $[a, b]$.

Proof of claim: Let $\varepsilon > 0$ be given. Let $t \in [a, b]$. Pick $t_m \in D \cap [a, b]$ such that in 24.47 $C_\varepsilon R |t - t_m| < \varepsilon/3$. Then there exists N such that if $l, n > N$, then $\|u_l(t_m) - u_n(t_m)\|_X < \varepsilon/3$. It follows that for $l, n > N$,

$$\begin{aligned} \|u_l(t) - u_n(t)\|_W &\leq \|u_l(t) - u_l(t_m)\|_W + \|u_l(t_m) - u_n(t_m)\|_W + \|u_n(t_m) - u_n(t)\|_W \\ &\leq \frac{2\varepsilon}{3} + \frac{\varepsilon}{3} + \frac{2\varepsilon}{3} < 2\varepsilon \end{aligned}$$

Since ε was arbitrary, this shows $\{u_k(t)\}_{k=1}^\infty$ is a Cauchy sequence. Since W is complete, this shows this sequence converges.

Now for $t \in [a, b]$, it was just shown that if $\varepsilon > 0$ there exists N_t such that if $n, m > N_t$, then $\|u_n(t) - u_m(t)\|_W < \frac{\varepsilon}{3}$. Now let $s \neq t$. Then

$$\|u_n(s) - u_m(s)\|_W \leq \|u_n(s) - u_n(t)\|_W + \|u_n(t) - u_m(t)\|_W + \|u_m(t) - u_m(s)\|_W$$

From 24.47

$$\|u_n(s) - u_m(s)\|_W \leq 2 \left(\frac{\varepsilon}{3} + C_\varepsilon R |t - s|^{1/q} \right) + \|u_n(t) - u_m(t)\|_W$$

and so it follows that if δ is sufficiently small and $s \in B(t, \delta)$, then when $n, m > N_t$ it follows that $\|u_n(s) - u_m(s)\|_W < \varepsilon$. Since $[a, b]$ is compact, there are finitely many of these balls, $\{B(t_i, \delta)\}_{i=1}^p$, such that for $s \in B(t_i, \delta)$ and $n, m > N_{t_i}$, the above inequality holds. Let $N > \max\{N_{t_1}, \dots, N_{t_p}\}$. Then if $m, n > N$ and $s \in [a, b]$ is arbitrary, it follows the above inequality must hold. Therefore, this has shown the following claim.

Claim: Let $\varepsilon > 0$ be given. Then there exists N such that if $m, n > N$, then it follows that $\|u_n - u_m\|_{\infty, W} < \varepsilon$.

Now let $u(t) = \lim_{k \rightarrow \infty} u_k(t)$.

$$\|u(t) - u(s)\|_W \leq \|u(t) - u_n(t)\|_W + \|u_n(t) - u_n(s)\|_W + \|u_n(s) - u(s)\|_W \quad (24.48)$$

Let N be in the above claim and fix $n > N$. Then

$$\|u(t) - u_n(t)\|_W = \lim_{m \rightarrow \infty} \|u_m(t) - u_n(t)\|_W \leq \varepsilon$$

and similarly, $\|u_n(s) - u(s)\|_W \leq \varepsilon$. Then if $|t - s|$ is small enough, 24.47 shows the middle term in 24.48 is also smaller than ε . Therefore, if $|t - s|$ is small enough,

$$\|u(t) - u(s)\|_W < 3\varepsilon.$$

Thus u is continuous. Finally, let N be as in the above claim. Then letting $m, n > N$, it follows that for all $t \in [a, b]$,

$$\|u_m(t) - u_n(t)\|_W < \varepsilon.$$

Therefore, letting $m \rightarrow \infty$, it follows that for all $t \in [a, b]$, $\|u(t) - u_n(t)\|_W \leq \varepsilon$. and so $\|u - u_n\|_{\infty, W} \leq \varepsilon$. ■

Here is an interesting corollary. Recall that for E a Banach space $C^{0, \alpha}([0, T], E)$ is the space of continuous functions u from $[0, T]$ to E such that $\|u\|_{\alpha, E} \equiv \|u\|_{\infty, E} + \rho_{\alpha, E}(u) < \infty$ where here $\rho_{\alpha, E}(u) \equiv \sup_{t \neq s} \frac{\|u(t) - u(s)\|_E}{|t - s|^\alpha}$

Corollary 24.11.7 *Let $E \subseteq W \subseteq X$ where the injection map is continuous from W to X and compact from E to W . Then if $\gamma > \alpha$, the embedding of $C^{0, \gamma}([0, T], E)$ into $C^{0, \alpha}([0, T], X)$ is compact.*

Proof: Let $\phi \in C^{0, \gamma}([0, T], E)$

$$\begin{aligned} \frac{\|\phi(t) - \phi(s)\|_X}{|t - s|^\alpha} &\leq \left(\frac{\|\phi(t) - \phi(s)\|_W}{|t - s|^\gamma} \right)^{\alpha/\gamma} \|\phi(t) - \phi(s)\|_W^{1 - (\alpha/\gamma)} \\ &\leq \left(\frac{\|\phi(t) - \phi(s)\|_E}{|t - s|^\gamma} \right)^{\alpha/\gamma} \|\phi(t) - \phi(s)\|_W^{1 - (\alpha/\gamma)} \leq \rho_{\gamma, E}(\phi) \|\phi(t) - \phi(s)\|_W^{1 - (\alpha/\gamma)} \end{aligned}$$

Now suppose $\{u_n\}$ is a bounded sequence in $C^{0, \gamma}([0, T], E)$. By Theorem 24.11.6 above, there is a subsequence still called $\{u_n\}$ which converges in $C^0([0, T], W)$. Thus from the above inequality

$$\begin{aligned} &\frac{\|u_n(t) - u_m(t) - (u_n(s) - u_m(s))\|_X}{|t - s|^\alpha} \\ &\leq \rho_{\gamma, E}(u_n - u_m) \|u_n(t) - u_m(t) - (u_n(s) - u_m(s))\|_W^{1 - (\alpha/\gamma)} \\ &\leq C(\{u_n\}) \left(2 \|u_n - u_m\|_{\infty, W} \right)^{1 - (\alpha/\gamma)} \end{aligned}$$

which converges to 0 as $n, m \rightarrow \infty$. Thus, $\rho_{\alpha, X}(u_n - u_m) \rightarrow 0$ as $n, m \rightarrow \infty$ Also

$$\|u_n - u_m\|_{\infty, X} \rightarrow 0$$

as $n, m \rightarrow \infty$ so this sequence is a Cauchy sequence in $C^{0, \alpha}([0, T], X)$. ■

The next theorem is a well known result probably due to Lions, Teman, or Aubin.

Theorem 24.11.8 *Let $E \subseteq W \subseteq X$ where the injection map is continuous from W to X and compact from E to W . Let $p \geq 1$, let $q > 1$, and define*

$$S \equiv \{u \in L^p([a, b]; E) : \text{for some } C, \|u(t) - u(s)\|_X \leq C|t - s|^{1/q} \text{ and } \|u\|_{L^p([a, b]; E)} \leq R\}.$$

Thus S is bounded in $L^p([a, b]; E)$ and Holder continuous into X . Then S is precompact in $L^p([a, b]; W)$. This means that if $\{u_n\}_{n=1}^\infty \subseteq S$, it has a subsequence $\{u_{n_k}\}$ which converges in $L^p([a, b]; W)$.

Proof: It suffices to show S has an η net in $L^p([a, b]; W)$ for each $\eta > 0$.

If not, there exists $\eta > 0$ and a sequence $\{u_n\} \subseteq S$, such that

$$\|u_n - u_m\| \geq \eta \quad (24.49)$$

for all $n \neq m$ and the norm refers to $L^p([a, b]; W)$. Let

$$a = t_0 < t_1 < \cdots < t_k = b, \quad t_i - t_{i-1} = (b - a)/k.$$

Now define $\bar{u}_n(t) \equiv \sum_{i=1}^k \bar{u}_{n_i} \mathcal{X}_{[t_{i-1}, t_i)}(t)$, $\bar{u}_{n_i} \equiv \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} u_n(s) ds$. The idea is to show that \bar{u}_n approximates u_n well and then to argue that a subsequence of the $\{\bar{u}_n\}$ is a Cauchy sequence yielding a contradiction to the above $\|u_n - u_m\| \geq \eta$.

Therefore,

$$\begin{aligned} u_n(t) - \bar{u}_n(t) &= \sum_{i=1}^k u_n(t) \mathcal{X}_{[t_{i-1}, t_i)}(t) - \sum_{i=1}^k \bar{u}_{n_i} \mathcal{X}_{[t_{i-1}, t_i)}(t) \\ &= \sum_{i=1}^k \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} u_n(t) ds \mathcal{X}_{[t_{i-1}, t_i)}(t) - \sum_{i=1}^k \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} u_n(s) ds \mathcal{X}_{[t_{i-1}, t_i)}(t) \\ &= \sum_{i=1}^k \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} (u_n(t) - u_n(s)) ds \mathcal{X}_{[t_{i-1}, t_i)}(t). \end{aligned}$$

It follows from Jensen's inequality, Lemma 10.15.1 that

$$\begin{aligned} \|u_n(t) - \bar{u}_n(t)\|_W^p &= \sum_{i=1}^k \left\| \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} (u_n(t) - u_n(s)) ds \right\|_W^p \mathcal{X}_{[t_{i-1}, t_i)}(t) \\ &\leq \sum_{i=1}^k \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} \|u_n(t) - u_n(s)\|_W^p ds \mathcal{X}_{[t_{i-1}, t_i)}(t) \end{aligned}$$

and so

$$\begin{aligned} &\int_a^b \|u_n(t) - \bar{u}_n(s)\|_W^p ds \\ &\leq \int_a^b \sum_{i=1}^k \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} \|u_n(t) - u_n(s)\|_W^p ds \mathcal{X}_{[t_{i-1}, t_i)}(t) dt \\ &= \sum_{i=1}^k \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} \int_{t_{i-1}}^{t_i} \|u_n(t) - u_n(s)\|_W^p ds dt. \end{aligned} \quad (24.50)$$

From Lemma 24.11.1 if $\varepsilon > 0$, there exists C_ε such that

$$\begin{aligned} \|u_n(t) - u_n(s)\|_W^p &\leq \varepsilon \|u_n(t) - u_n(s)\|_E^p + C_\varepsilon \|u_n(t) - u_n(s)\|_X^p \\ &\leq 2^{p-1} \varepsilon (\|u_n(t)\|^p + \|u_n(s)\|^p) + C_\varepsilon |t - s|^{p/q} \end{aligned}$$

This is substituted in to 24.50 to obtain

$$\int_a^b \|u_n(t) - \bar{u}_n(s)\|_W^p ds \leq$$

$$\begin{aligned}
& \sum_{i=1}^k \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} \int_{t_{i-1}}^{t_i} \left(2^{p-1} \varepsilon (\|u_n(t)\|^p + \|u_n(s)\|^p) + C_\varepsilon |t-s|^{p/q} \right) ds dt \\
&= \sum_{i=1}^k 2^p \varepsilon \int_{t_{i-1}}^{t_i} \|u_n(t)\|_W^p dt + \frac{C_\varepsilon}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} \int_{t_{i-1}}^{t_i} |t-s|^{p/q} ds dt \\
&\leq 2^p \varepsilon \int_a^b \|u_n(t)\|^p dt + C_\varepsilon \sum_{i=1}^k \frac{1}{(t_i - t_{i-1})} (t_i - t_{i-1})^{p/q} \int_{t_{i-1}}^{t_i} \int_{t_{i-1}}^{t_i} ds dt \\
&= 2^p \varepsilon \int_a^b \|u_n(t)\|^p dt + C_\varepsilon \sum_{i=1}^k \frac{1}{(t_i - t_{i-1})} (t_i - t_{i-1})^{p/q} (t_i - t_{i-1})^2 \\
&\leq 2^p \varepsilon R^p + C_\varepsilon \sum_{i=1}^k (t_i - t_{i-1})^{1+p/q} = 2^p \varepsilon R^p + C_\varepsilon k \left(\frac{b-a}{k} \right)^{1+p/q}.
\end{aligned}$$

Taking ε so small that $2^p \varepsilon R^p < \eta^p/8^p$ and then choosing k sufficiently large, it follows that $\|u_n - \bar{u}_n\|_{L^p([a,b];W)} < \frac{\eta}{4}$.

Thus k is fixed and \bar{u}_n at a step function with k steps having values in E . Now use compactness of the embedding of E into W to obtain a subsequence such that $\{\bar{u}_n\}$ is Cauchy in $L^p(a, b; W)$ and use this to contradict 24.49. The details follow.

Suppose $\bar{u}_n(t) = \sum_{i=1}^k u_i^n \mathcal{X}_{[t_{i-1}, t_i)}(t)$. Thus $\|\bar{u}_n(t)\|_E = \sum_{i=1}^k \|u_i^n\|_E \mathcal{X}_{[t_{i-1}, t_i)}(t)$ and so $R \geq \int_a^b \|\bar{u}_n(t)\|_E^p dt = \frac{T}{k} \sum_{i=1}^k \|u_i^n\|_E^p$. Therefore, the $\{u_i^n\}$ are all bounded. It follows that after taking subsequences k times there exists a subsequence $\{u_{n_k}\}$ such that u_{n_k} is a Cauchy sequence in $L^p(a, b; W)$. You simply get a subsequence such that $u_i^{n_k}$ is a Cauchy sequence in W for each i . Then denoting this subsequence by n ,

$$\begin{aligned}
\|u_n - u_m\|_{L^p(a,b;W)} &\leq \|u_n - \bar{u}_n\|_{L^p(a,b;W)} + \|\bar{u}_n - \bar{u}_m\|_{L^p(a,b;W)} + \|\bar{u}_m - u_m\|_{L^p(a,b;W)} \\
&\leq \frac{\eta}{4} + \|\bar{u}_n - \bar{u}_m\|_{L^p(a,b;W)} + \frac{\eta}{4} < \eta
\end{aligned}$$

provided m, n are large enough, contradicting 24.49. ■

You can give a different version of the above to include the case where there is, instead of a Holder condition, a bound on u' for $u \in S$. It is stated next. We are assuming a situation in which $\int_a^b u'(t) dt = u(b) - u(a)$. This happens, for example, if u' is the weak derivative. This is discussed in the exercises. These kind of theorems are in [52].

Corollary 24.11.9 *Let $E \subseteq W \subseteq X$ where the injection map is continuous from W to X and compact from E to W . Let $p \geq 1$, let $q > 1$, and define*

$$S \equiv \{u \in L^p([a, b]; E) : \text{for some } C, \|u(t) - u(s)\|_X \leq C|t-s|^{1/q} \text{ and } \|u\|_{L^p([a,b];E)} \leq R\}.$$

Thus S is bounded in $L^p([a, b]; E)$ and Holder continuous into X . Then S is precompact in $L^p([a, b]; W)$. This means that if $\{u_n\}_{n=1}^\infty \subseteq S$, it has a subsequence $\{u_{n_k}\}$ which converges in $L^p([a, b]; W)$. The same conclusion can be drawn if it is known instead of the Holder condition that $\|u'\|_{L^1([a,b];X)}$ is bounded.

Proof: The first part is Theorem 24.11.8. Therefore, we just prove the new stuff which involves a bound on the L^1 norm of the derivative. It suffices to show S has an η net in $L^p([a, b]; W)$ for each $\eta > 0$.

If not, there exists $\eta > 0$ and a sequence $\{u_n\} \subseteq S$, such that

$$\|u_n - u_m\| \geq \eta \quad (24.51)$$

for all $n \neq m$ and the norm refers to $L^p([a, b]; W)$. Let

$$a = t_0 < t_1 < \cdots < t_k = b, \quad t_i - t_{i-1} = (b - a) / k.$$

Now define $\bar{u}_n(t) \equiv \sum_{i=1}^k \bar{u}_{n_i} \mathcal{X}_{[t_{i-1}, t_i)}(t)$, $\bar{u}_{n_i} \equiv \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} u_n(s) ds$. The idea is to show that \bar{u}_n approximates u_n well and then to argue that a subsequence of the $\{\bar{u}_n\}$ is a Cauchy sequence yielding a contradiction to 24.51.

Therefore,

$$\begin{aligned} u_n(t) - \bar{u}_n(t) &= \sum_{i=1}^k u_n(t) \mathcal{X}_{[t_{i-1}, t_i)}(t) - \sum_{i=1}^k \bar{u}_{n_i} \mathcal{X}_{[t_{i-1}, t_i)}(t) \\ &= \sum_{i=1}^k \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} u_n(t) ds \mathcal{X}_{[t_{i-1}, t_i)}(t) - \sum_{i=1}^k \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} u_n(s) ds \mathcal{X}_{[t_{i-1}, t_i)}(t) \\ &= \sum_{i=1}^k \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} (u_n(t) - u_n(s)) ds \mathcal{X}_{[t_{i-1}, t_i)}(t). \end{aligned}$$

It follows from Jensen's inequality, Lemma 10.15.1, that

$$\|u_n(t) - \bar{u}_n(t)\|_W^p = \sum_{i=1}^k \left\| \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} (u_n(t) - u_n(s)) ds \right\|_W^p \mathcal{X}_{[t_{i-1}, t_i)}(t)$$

And so

$$\begin{aligned} \int_0^T \|u_n(t) - \bar{u}_n(t)\|_W^p dt &= \sum_{i=1}^k \int_{t_{i-1}}^{t_i} \left\| \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} (u_n(t) - u_n(s)) ds \right\|_W^p dt \\ &\leq \sum_{i=1}^k \int_{t_{i-1}}^{t_i} \varepsilon \left\| \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} (u_n(t) - u_n(s)) ds \right\|_E^p dt \\ &\quad + C_\varepsilon \sum_{i=1}^k \int_{t_{i-1}}^{t_i} \left\| \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} (u_n(t) - u_n(s)) ds \right\|_X^p dt \end{aligned} \quad (24.52)$$

Consider the second of these. It equals $C_\varepsilon \sum_{i=1}^k \int_{t_{i-1}}^{t_i} \left\| \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} \int_s^t u'_n(\tau) d\tau ds \right\|_X^p dt$. This is no larger than

$$\begin{aligned} &\leq C_\varepsilon \sum_{i=1}^k \int_{t_{i-1}}^{t_i} \left(\frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} \int_{t_{i-1}}^{t_i} \|u'_n(\tau)\|_X d\tau ds \right)^p dt \\ &= C_\varepsilon \sum_{i=1}^k \int_{t_{i-1}}^{t_i} \left(\int_{t_{i-1}}^{t_i} \|u'_n(\tau)\|_X d\tau \right)^p dt = C_\varepsilon \sum_{i=1}^k \left((t_i - t_{i-1})^{1/p} \int_{t_{i-1}}^{t_i} \|u'_n(\tau)\|_X d\tau \right)^p \end{aligned}$$

Since $\frac{b-a}{k} = t_i - t_{i-1}$,

$$\begin{aligned} &= C_\varepsilon \left(\sum_{i=1}^k \left(\frac{b-a}{k} \right)^{1/p} \int_{t_{i-1}}^{t_i} \|u'_n(\tau)\|_X d\tau \right)^p \leq \frac{C_\varepsilon (b-a)}{k} \left(\sum_{i=1}^k \int_{t_{i-1}}^{t_i} \|u'_n(\tau)\|_X d\tau \right)^p \\ &= \frac{C_\varepsilon (b-a)}{k} \left(\|u'_n\|_{L^1([a, b], X)} \right)^p < \frac{\eta^p}{10^p} \end{aligned}$$

if k is chosen large enough. Now consider the first in 24.52. By Jensen's inequality, Lemma 10.15.1,

$$\begin{aligned}
& \sum_{i=1}^k \int_{t_{i-1}}^{t_i} \varepsilon \left\| \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} (u_n(t) - u_n(s)) ds \right\|_E^p dt \\
& \leq \sum_{i=1}^k \int_{t_{i-1}}^{t_i} \varepsilon \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} \|u_n(t) - u_n(s)\|_E^p ds dt \\
& \leq \varepsilon 2^{p-1} \sum_{i=1}^k \frac{1}{t_i - t_{i-1}} \int_{t_{i-1}}^{t_i} \int_{t_{i-1}}^{t_i} (\|u_n(t)\|^p + \|u_n(s)\|^p) ds dt \\
& = 2\varepsilon 2^{p-1} \sum_{i=1}^k \int_{t_{i-1}}^{t_i} (\|u_n(t)\|^p) dt = \varepsilon (2) (2^{p-1}) \|u_n\|_{L^p([a,b],E)}^p \leq M\varepsilon
\end{aligned}$$

Now pick ε sufficiently small that $M\varepsilon < \frac{\eta^p}{10^p}$ and then k large enough that the second term in 24.52 is also less than $\eta^p/10^p$. Then it will follow that

$$\|\bar{u}_n - u_n\|_{L^p([a,b],W)} < \left(\frac{2\eta^p}{10^p} \right)^{1/p} = 2^{1/p} \frac{\eta}{10} \leq \frac{\eta}{5}$$

Thus k is fixed and \bar{u}_n at a step function with k steps having values in E . Now use compactness of the embedding of E into W to obtain a subsequence such that $\{\bar{u}_n\}$ is Cauchy in $L^p([a,b];W)$ and use this to contradict 24.51. The details follow.

Suppose $\bar{u}_n(t) = \sum_{i=1}^k u_i^n \mathcal{X}_{[t_{i-1}, t_i)}(t)$. Thus $\|\bar{u}_n(t)\|_E = \sum_{i=1}^k \|u_i^n\|_E \mathcal{X}_{[t_{i-1}, t_i)}(t)$ and so $R \geq \int_a^b \|\bar{u}_n(t)\|_E^p dt = \frac{T}{k} \sum_{i=1}^k \|u_i^n\|_E^p$. Therefore, the $\{u_i^n\}$ are all bounded. It follows that after taking subsequences k times there exists a subsequence $\{u_{n_k}\}$ such that u_{n_k} is a Cauchy sequence in $L^p([a,b];W)$. You simply get a subsequence such that $u_i^{n_k}$ is a Cauchy sequence in W for each i . Then denoting this subsequence by n ,

$$\begin{aligned}
\|u_n - u_m\|_{L^p(a,b;W)} & \leq \|u_n - \bar{u}_n\|_{L^p(a,b;W)} + \|\bar{u}_n - \bar{u}_m\|_{L^p(a,b;W)} + \|\bar{u}_m - u_m\|_{L^p(a,b;W)} \\
& \leq \frac{\eta}{4} + \|\bar{u}_n - \bar{u}_m\|_{L^p(a,b;W)} + \frac{\eta}{4} < \eta
\end{aligned}$$

provided m, n are large enough, contradicting 24.51. ■

24.12 Conditional Expectation in Banach Spaces

Let (Ω, \mathcal{F}, P) be a probability space and let $X \in L^1(\Omega; \mathbb{R})$. Also let $\mathcal{G} \subseteq \mathcal{F}$ where \mathcal{G} is also a σ algebra. Then the usual conditional expectation is defined by

$$\int_A X dP = \int_A E(X|\mathcal{G}) dP$$

where $E(X|\mathcal{G})$ is \mathcal{G} measurable and $A \in \mathcal{G}$ is arbitrary. Recall this is an application of the Radon Nikodym theorem. Also recall $E(X|\mathcal{G})$ is unique up to a set of measure zero.

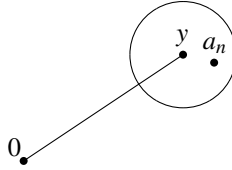
I want to do something like this here. Denote by $L^1(\Omega; E, \mathcal{G})$ those functions in $L^1(\Omega; E)$ which are measurable with respect to \mathcal{G} .

Theorem 24.12.1 Let E be a separable Banach space and let $X \in L^1(\Omega; E, \mathcal{F})$ where X is measurable with respect to \mathcal{F} and let \mathcal{G} be a σ algebra which is contained in \mathcal{F} . Then there exists a unique $Z \in L^1(\Omega; E, \mathcal{G})$ such that for all $A \in \mathcal{G}$,

$$\int_A X dP = \int_A Z dP$$

Denoting this Z as $E(X|\mathcal{G})$, it follows $\|E(X|\mathcal{G})\| \leq E(\|X\| |\mathcal{G})$.

Proof: First consider uniqueness. Suppose Z' is another in $L^1(\Omega; E, \mathcal{G})$ which works. Consider a dense subset of E $\{a_n\}_{n=1}^\infty$. Then the balls $\left\{B\left(a_n, \frac{\|a_n\|}{4}\right)\right\}_{n=1}^\infty$ must cover $E \setminus \{0\}$. Here is why. If $y \neq 0$, pick $a_n \in B\left(y, \frac{\|y\|}{5}\right)$.



Then $\|a_n\| \geq 4\|y\|/5$ and so $\|a_n - y\| < \|y\|/5$. Thus $y \in B(a_n, \|y\|/5) \subseteq B\left(a_n, \frac{\|a_n\|}{4}\right)$. Now suppose Z is \mathcal{G} measurable and $\int_A Z dP = 0$ for all $A \in \mathcal{G}$. Then define the set A by $A \equiv Z^{-1}\left(B\left(a_n, \frac{\|a_n\|}{4}\right)\right)$ it follows $0 = \int_A Z - a_n + a_n dP$ and so

$$\begin{aligned} \|a_n\| P(A) &= \left\| \int_A a_n dP \right\| = \left\| \int_A (a_n - Z) dP \right\| \\ &\leq \int_{Z^{-1}\left(B\left(a_n, \frac{\|a_n\|}{4}\right)\right)} \|a_n - Z\| dP \leq \frac{\|a_n\|}{4} P(A) \end{aligned}$$

which is a contradiction unless $P(A) = 0$. Therefore, letting

$$N \equiv \bigcup_{n=1}^\infty Z^{-1}\left(B\left(a_n, \frac{\|a_n\|}{4}\right)\right) = Z^{-1}(E \setminus \{0\})$$

it follows N has measure zero and so $Z = 0$ a.e. This proves uniqueness because if Z, Z' both hold, then from the above argument, $Z - Z' = 0$ a.e.

Next I will show Z exists. To do this recall Theorem 24.2.4 on Page 656 which is stated below for convenience.

Theorem 24.12.2 An E valued function, X , is Bochner integrable if and only if X is strongly measurable and

$$\int_\Omega \|X(\omega)\| dP < \infty. \quad (24.53)$$

In this case there exists a sequence of simple functions $\{X_n\}$ satisfying

$$\int_\Omega \|X_n(\omega) - X_m(\omega)\| dP \rightarrow 0 \text{ as } m, n \rightarrow \infty. \quad (24.54)$$

$X_n(\omega)$ converging pointwise to $X(\omega)$,

$$\|X_n(\omega)\| \leq 2\|X(\omega)\| \quad (24.55)$$

and

$$\lim_{n \rightarrow \infty} \int_{\Omega} \|X(\omega) - X_n(\omega)\| dP = 0. \quad (24.56)$$

Now let $\{X_n\}$ be the simple functions just defined and let $X_n(\omega) = \sum_{k=1}^m x_k \mathcal{X}_{F_k}(\omega)$ where $F_k \in \mathcal{F}$, the F_k being disjoint. Then define $Z_n \equiv \sum_{k=1}^m x_k E(\mathcal{X}_{F_k} | \mathcal{G})$. Thus, if $A \in \mathcal{G}$,

$$\begin{aligned} \int_A Z_n dP &= \sum_{k=1}^m x_k \int_A E(\mathcal{X}_{F_k} | \mathcal{G}) dP = \sum_{k=1}^m x_k \int_A \mathcal{X}_{F_k} dP \\ &= \sum_{k=1}^m x_k P(F_k \cap A) = \int_A X_n dP \end{aligned} \quad (24.57)$$

Then since $E(\mathcal{X}_{F_k} | \mathcal{G}) \geq 0$, it follows that $\|Z_n\| \leq \sum_{k=1}^m \|x_k\| E(\mathcal{X}_{F_k} | \mathcal{G})$. Thus if $A \in \mathcal{G}$,

$$\begin{aligned} E(\|Z_n\| | \mathcal{A}) &\leq E\left(\sum_{k=1}^m \|x_k\| \mathcal{X}_A E(\mathcal{X}_{F_k} | \mathcal{G})\right) = \sum_{k=1}^m \|x_k\| \int_A E(\mathcal{X}_{F_k} | \mathcal{G}) dP \\ &= \sum_{k=1}^m \|x_k\| \int_A \mathcal{X}_{F_k} dP = E(\mathcal{X}_A \|X_n\|). \end{aligned} \quad (24.58)$$

Note the use of \leq in the first step in the above. Although the F_k are disjoint, all that is known about $E(\mathcal{X}_{F_k} | \mathcal{G})$ is that it is nonnegative. Similarly, $E(\|Z_n - Z_m\|) \leq E(\|X_n - X_m\|)$ and this last term converges to 0 as $n, m \rightarrow \infty$ by the properties of the X_n . Therefore, $\{Z_n\}$ is a Cauchy sequence in $L^1(\Omega; E; \mathcal{G})$. It follows it converges to some Z in $L^1(\Omega; E; \mathcal{G})$. Then letting $A \in \mathcal{G}$, and using 24.57,

$$\begin{aligned} \int_A Z dP &= \int \mathcal{X}_A Z dP = \lim_{n \rightarrow \infty} \int \mathcal{X}_A Z_n dP = \lim_{n \rightarrow \infty} \int_A Z_n dP \\ &= \lim_{n \rightarrow \infty} \int_A X_n dP = \int_A X dP. \end{aligned}$$

Then define $Z \equiv E(X | \mathcal{G})$.

It remains to verify $\|E(X | \mathcal{G})\| \equiv \|Z\| \leq E(\|X\| | \mathcal{G})$. This follows because, from the above, $\|Z_n\| \rightarrow \|Z\|$, $\|X_n\| \rightarrow \|X\|$ in $L^1(\Omega)$ and so if $A \in \mathcal{G}$, then from 24.58,

$$\frac{1}{P(A)} \int_A \|Z_n\| dP \leq \frac{1}{P(A)} \int_A \|X_n\| dP$$

and so, passing to the limit,

$$\frac{1}{P(A)} \int_A \|Z\| dP \leq \frac{1}{P(A)} \int_A \|X\| dP = \frac{1}{P(A)} \int_A E(\|X\| | \mathcal{G}) dP$$

Since A is arbitrary, this shows that $\|E(X | \mathcal{G})\| \equiv \|Z\| \leq E(\|X\| | \mathcal{G})$. ■

In the case where E is reflexive, one could also use Corollary 24.7.6 on Page 682 to get the above result. You would define a vector measure on \mathcal{G} , $\nu(F) \equiv \int_F X dP$ and then you would use the fact that reflexive separable Banach spaces have the Radon Nikodym property to obtain $Z \in L^1(\Omega; E; \mathcal{G})$ such that $\nu(F) = \int_F X dP = \int_F Z dP$.

The function, Z whose existence and uniqueness is guaranteed by Theorem 24.12.2 is called $E(X | \mathcal{G})$.

24.13 Exercises

1. Show $L^1(\mathbb{R})$ is not reflexive. **Hint:** $L^1(\mathbb{R})$ is separable. What about $L^\infty(\mathbb{R})$?
2. If $f \in L^1(\mathbb{R}^n; X)$ for X a Banach space, does the usual fundamental theorem of calculus work? That is, can you say $\lim_{r \rightarrow 0} \frac{1}{m(B(x, r))} \int_{B(x, r)} f(t) dm = f(x)$ a.e.?
3. Does the Vitali convergence theorem hold for Bochner integrable functions? If so, give a statement of the appropriate theorem and a proof.
4. Suppose $g \in L^1([a, b]; X)$ where X is a Banach space. Then if $\int_a^b g(t) \phi(t) dt = 0$ for all $\phi \in C_c^\infty(a, b)$, then $g(t) = 0$ a.e. Show that this is the case. **Hint:** It will likely depend on the regularity properties of Lebesgue measure.
5. Suppose $f \in L^1(a, b; X)$ and for all $\phi \in C_c^\infty(a, b)$, $\int_a^b f(t) \phi'(t) dt = 0$. Then there exists a constant, $a \in X$ such that $f(t) = a$ a.e. **Hint:** Let

$$\psi_\phi(x) \equiv \int_a^x [\phi(t) - \left(\int_a^b \phi(y) dy \right) \phi_0(t)] dt, \quad \phi_0 \in C_c^\infty(a, b), \quad \int_a^b \phi_0(x) dx = 1$$

Then explain why $\psi_\phi \in C_c^\infty(a, b)$, $\psi'_\phi = \phi - \left(\int_a^b \phi(y) dy \right) \phi_0$. Then use the assumption on ψ_ϕ . Next use the above problem. Verify that $f(y) = \int_a^b f(t) \phi_0(t) dt$ a.e. y

6. Let $f \in L^1([a, b], X)$. Then we say that the weak derivative of f is in $L^1([a, b], X)$ if there is a function denoted as $f' \in L^1([a, b], X)$ such that for all $\phi \in C_c^\infty(a, b)$,

$$-\int_a^b f(t) \phi'(t) dt = \int_a^b f'(t) \phi(t) dt$$

Show that this definition is well defined. Next, using the above problems, show that if $f, f' \in L^1([a, b], X)$, it follows that there is a continuous function, denoted by $t \rightarrow \hat{f}(t)$ such that $\hat{f}(t) = f(t)$ a.e. t and $\hat{f}(t) = \hat{f}(a) + \int_0^t f'(s) ds$. Thus, unlike the classical definition of the derivative, when a function and its derivative are both in L^1 , it has a representative \hat{f} which equals the function a.e. such that \hat{f} can be recovered from its derivative. Recall the well known example of this not working out which is based on the Cantor function of Problem 4 on Page 268. This function had zero derivative a.e. and yet it climbed from 0 to 1 on the unit interval. Thus one could not recover it from integrating its classical derivative. Incidentally, if the function has a derivative everywhere, then you can recover it by taking the generalized Riemann integral of the derivative, although the Lebesgue integral of this derivative might not even be defined. This is in my book on single variable advanced calculus, but this integral is not discussed here.

Chapter 25

Stone's Theorem and Partitions of Unity

This section is devoted to Stone's theorem which says that a metric space is paracompact, defined below. See [41] for this which is where I read it. First is the definition of what is meant by a refinement. A metric space is an example of a topological space and it is the context for what is done below.

Definition 25.0.1 *Let S be a topological space. We say that a collection of sets \mathfrak{D} is a refinement of an open cover \mathfrak{S} , if every set of \mathfrak{D} is contained in some set of \mathfrak{S} . An open refinement would be one in which all sets are open, with a similar convention holding for the term "closed refinement".*

Definition 25.0.2 *We say that a collection of sets \mathfrak{D} , is locally finite if for all $p \in S$, the topological space, there exists V an open set containing p such that V has nonempty intersection with only finitely many sets of \mathfrak{D} .*

Definition 25.0.3 *We say S is paracompact if it is Hausdorff and for every open cover \mathfrak{S} , there exists an open refinement \mathfrak{D} such that \mathfrak{D} is locally finite and \mathfrak{D} covers S .*

Recall how the union of finitely many closed sets is closed. This can be generalized to a locally finite set of closed sets. Think \mathbb{N} for example. The following implies this.

Theorem 25.0.4 *If \mathfrak{D} is locally finite then*

$$\cup\{\overline{D} : D \in \mathfrak{D}\} = \overline{\cup\{D : D \in \mathfrak{D}\}}.$$

Proof: It is clear the left side is a subset of the right because the right is a closed set which contains the left since a limit point of any D is in the set on the right. If $p \in \cup\{D : D \in \mathfrak{D}\}$, there is nothing to show. Let p not be in this set but be a limit point of $\cup\{D : D \in \mathfrak{D}\}$. Is p in some \overline{D} ? Let $p \in V$, an open set intersecting only finitely many sets of \mathfrak{D} , D_1, \dots, D_n . If p is not in any of \overline{D}_i then $p \in W$ where W is some open set which contains no points of $\cup_{i=1}^n D_i$. Then $V \cap W$ contains no points of any set of \mathfrak{D} and this contradicts the assumption that p is a limit point of $\cup\{D : D \in \mathfrak{D}\}$. Thus $p \in \overline{D}_i$ for some i . ■

We say $\mathfrak{S} \subseteq \mathcal{P}(S)$ is countably locally finite if

$$\mathfrak{S} = \cup_{n=1}^{\infty} \mathfrak{S}_n$$

and each \mathfrak{S}_n is locally finite. The following theorem appeared in the 1950's. It will be used to prove Stone's theorem.

Theorem 25.0.5 *Let S be a regular topological space. (If $p \in U$ open, then there exists an open set V such that $p \in \bar{V} \subseteq U$.) The following are equivalent*

- 1.) *Every open covering of S has a refinement that is open, covers S and is **countably** locally finite.*
- 2.) *Every open covering of S has a refinement that is locally finite and covers S . (The sets in refinement maybe not open.)*
- 3.) *Every open covering of S has a refinement that is closed, locally finite, and covers S . (Sets in refinement are closed.)*
- 4.) *Every open covering of S has a refinement that is open, locally finite, and covers S . (Sets in refinement are open.)*

Proof:

1.) \Rightarrow 2.)

Let \mathfrak{S} be an open cover of S and let \mathfrak{B} be an open countably locally finite refinement

$$\mathfrak{B} = \bigcup_{n=1}^{\infty} \mathfrak{B}_n$$

where \mathfrak{B}_n is an open refinement of \mathfrak{S} and \mathfrak{B}_n is locally finite. For $B \in \mathfrak{B}_n$, let

$$E_n(B) = B \setminus \bigcup_{k < n} (\bigcup \{B : B \in \mathfrak{B}_k\}) \equiv B \cap \left(\bigcup_{k < n} (\bigcup \{B : B \in \mathfrak{B}_k\}) \right)^C.$$

Thus, in words, $E_n(B)$ consists of points in B which are not in any set from any \mathfrak{B}_k for $k < n$.

Claim: $\{E_n(B) : n \in \mathbb{N}, B \in \mathfrak{B}_n\}$ is locally finite.

Proof of the claim: Let $p \in S$. Then $p \in B_0 \in \mathfrak{B}_n$ for some n . Let V be open, $p \in V$, and V intersects only finitely many sets of $\mathfrak{B}_1 \cup \dots \cup \mathfrak{B}_n$. Then consider $B_0 \cap V$. If $m > n$,

$$(B_0 \cap V) \cap E_m(B) \subseteq \left(\bigcup_{k < m} (\bigcup \{B : B \in \mathfrak{B}_k\}) \right)^C \subseteq B_0^C.$$

In words, $E_m(B)$ has nothing in it from any of the \mathfrak{B}_k for $k < m$. In particular, it has nothing in it from B_0 . Thus $(B_0 \cap V) \cap E_m(B) = \emptyset$ for $m > n$. Thus $p \in B_0 \cap V$ which intersects only finitely many sets of \mathfrak{S} , no more than those intersected by V . This establishes the claim.

Claim: $\{E_n(B) : n \in \mathbb{N}, B \in \mathfrak{B}_n\}$ covers S .

Proof: Let $p \in S$ and let $n = \min\{k \in \mathbb{N} : p \in B \text{ for some } B \in \mathfrak{B}_k\}$. Let $p \in B \in \mathfrak{B}_n$. Then $p \in E_n(B)$.

The two claims show that 1.) \Rightarrow 2.).

2.) \Rightarrow 3.)

Let \mathfrak{S} be an open cover and let

$$\mathcal{G} \equiv \{U : U \text{ is open and } \overline{U} \subseteq V \in \mathfrak{S} \text{ for some } V \in \mathfrak{S}\}.$$

Then since S is regular, \mathcal{G} covers S . (If $p \in S$, then $p \in U \subseteq \overline{U} \subseteq V \in \mathfrak{S}$.) By 2.), \mathcal{G} has a locally finite refinement \mathfrak{C} , covering S . Consider $\{\overline{E} : E \in \mathfrak{C}\}$. This collection of closed sets covers S and is locally finite because if $p \in S$, there exists V , $p \in V$, and V has nonempty intersections with only finitely many elements of \mathfrak{C} , say E_1, \dots, E_n . If $\overline{E} \cap V \neq \emptyset$, then $E \cap V \neq \emptyset$ and so V intersects only $\overline{E}_1, \dots, \overline{E}_n$. This shows 2.) \Rightarrow 3.).

3.) \Rightarrow 4.) Here is a table of symbols with a short summary of their meaning.

Open covering	Locally finite refinement
\mathfrak{S} original covering	\mathfrak{B} by 3. can be closed refinement
\mathfrak{F} open intersectors	\mathfrak{C} closed refinement

Let \mathfrak{S} be an open cover and let \mathfrak{B} be a locally finite refinement which covers S . By 3.) we can take \mathfrak{B} to be a closed refinement but this is not important here. Let

$$\mathfrak{F} \equiv \{U : U \text{ is open and } U \text{ intersects only finitely many sets of } \mathfrak{B}\}.$$

Then \mathfrak{F} covers S because \mathfrak{B} is locally finite. (If $p \in S$, then there exists an open set U containing p which intersects only finitely many sets of \mathfrak{B} . Thus $p \in U \in \mathfrak{F}$.) By 3., \mathfrak{F} has a locally finite closed refinement \mathfrak{C} , which covers S . Define for $B \in \mathfrak{B}$

$$\mathfrak{C}(B) \equiv \{C \in \mathfrak{C} : C \cap B = \emptyset\}$$

Thus these closed sets C do not intersect B and so B is in their complement. We use $\mathfrak{C}(B)$ to fatten up B . Let

$$E(B) \equiv (\cup \{C : C \in \mathfrak{C}(B)\})^C.$$

In words, $E(B)$ is the complement of the union of all closed sets of \mathfrak{C} which do not intersect B . Thus $E(B) \supseteq B$, and has fattened up B . Then since $\mathfrak{C}(B)$ is locally finite, $E(B)$ is an open set by Theorem 25.0.4. Now let $F(B)$ be defined such that for $B \in \mathfrak{B}$,

$$B \subseteq F(B) \in \mathfrak{S}$$

(by definition B is in some set of \mathfrak{S}), and let

$$\mathfrak{L} = \{E(B) \cap F(B) : B \in \mathfrak{B}\}$$

The intersection with $F(B)$ is to ensure that \mathfrak{L} is a refinement of \mathfrak{S} . The important thing to notice is that **if $C \in \mathfrak{C}$ intersects $E(B)$, then it must also intersect B** . If not, you could include it in the list of closed sets which do not intersect B and whose complement is $E(B)$. Thus $E(B)$ would be too large.

Claim: \mathfrak{L} covers S .

This claim is obvious because if $p \in S$ then $p \in B$ for some $B \in \mathfrak{B}$. Hence

$$p \in E(B) \cap F(B) \in \mathfrak{L}.$$

Claim: \mathfrak{L} is locally finite and a refinement of \mathfrak{S} .

Proof: It is clear \mathfrak{L} is a refinement of \mathfrak{S} because every set of \mathfrak{L} is a subset of a set of \mathfrak{S} , $F(B)$. Let $p \in S$. There exists an open set W , such that $p \in W$ and W intersects only C_1, \dots, C_n , elements of \mathfrak{C} . Hence $W \subseteq \cup_{i=1}^n C_i$ since \mathfrak{C} covers S .

But C_i is contained in a set $U_i \in \mathfrak{F}$ which intersects only finitely many sets of \mathfrak{B} . Thus each C_i intersects only finitely many $B \in \mathfrak{B}$ and so each C_i intersects only finitely many of the sets, $E(B)$. (If it intersects $E(B)$, then it intersects B .) Thus W intersects only finitely many of the $E(B)$, hence finitely many of the $E(B) \cap F(B)$. It follows that \mathfrak{L} is locally finite.

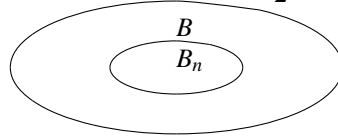
It is obvious that 4.) \Rightarrow 1.). ■

The following theorem is Stone's theorem.

Theorem 25.0.6 *If S is a metric space then S is paracompact (Every open cover has a locally finite open refinement also an open cover.)*

Proof: Let \mathfrak{S} be an open cover. Well order \mathfrak{S} . For $B \in \mathfrak{S}$,

$$B_n \equiv \{x \in B : \text{dist}(x, B^C) < \frac{1}{2^n}\}, n = 1, 2, \dots$$



Thus B_n is contained in B but approximates it up to 2^{-n} . Let

$$E_n(B) = B_n \setminus \cup\{D : D \prec B \text{ and } D \neq B\}$$

where \prec denotes the well order. If $B, D \in \mathfrak{S}$, then one is first in the well order. Let $D \prec B$. Then from the construction, $E_n(B) \subseteq D^C$ and $E_n(D)$ is further than $1/2^n$ from D^C . Hence, assuming neither set is empty,

$$\text{dist}(E_n(B), E_n(D)) \geq 2^{-n}$$

for all $B, D \in \mathfrak{S}$. Fatten up $E_n(B)$ as follows.

$$\widetilde{E_n(B)} \equiv \cup\{B(x, 8^{-n}) : x \in E_n(B)\}.$$

Thus $\widetilde{E_n(B)} \subseteq B$ and

$$\text{dist}(\widetilde{E_n(B)}, \widetilde{E_n(D)}) \geq \frac{1}{2^n} - 2\left(\frac{1}{8}\right)^n \equiv \delta_n > 0.$$

It follows that the collection of open sets

$$\{\widetilde{E_n(B)} : B \in \mathfrak{S}\} \equiv \mathfrak{B}_n$$

is locally finite. In fact, $B(p, \frac{\delta_n}{2})$ cannot intersect more than one of them. In addition to this,

$$S \subseteq \cup\{\widetilde{E_n(B)} : n \in \mathbb{N}, B \in \mathfrak{S}\}$$

because if $p \in S$, let B be the first set in \mathfrak{S} to contain p . Then $p \in E_n(B)$ for n large enough because it will not be in anything deleted. Thus this is an open countably locally finite refinement. Thus 1.) in the above theorem is satisfied. ■

25.1 Partitions of Unity and Stone's Theorem

First recall that if S is nonempty, then $\text{dist}(x, S)$ satisfies $|\text{dist}(x, S) - \text{dist}(y, S)| \leq d(x, y)$. It was Lemma 3.12.1.

Theorem 25.1.1 *Let S be a metric space and let \mathfrak{S} be any open cover of S . Then there exists a set \mathfrak{F} , an open refinement of \mathfrak{S} , and functions $\{\phi_F : F \in \mathfrak{F}\}$ such that*

$$\phi_F : S \rightarrow [0, 1]$$

$$\phi_F \text{ is continuous}$$

$$\phi_F(x) \text{ equals 0 for all but finitely many } F \in \mathfrak{F}$$

$$\sum\{\phi_F(x) : F \in \mathfrak{F}\} = 1 \text{ for all } x \in S.$$

Each ϕ_F is locally Lipschitz continuous which means that for each z there is an open set W containing z for which, if $x, y \in W$, then there is a constant K such that

$$|\phi_F(x) - \phi_F(y)| \leq Kd(x, y)$$

Proof: By Stone's theorem, there exists a locally finite open refinement \mathfrak{F} of \mathfrak{S} covering S . For $F \in \mathfrak{F}$

$$g_F(x) \equiv \text{dist}(x, F^C)$$

Let

$$\phi_F(x) \equiv \left(\sum \{g_F(x) : F \in \mathfrak{F}\} \right)^{-1} g_F(x).$$

Now

$$\sum \{g_F(x) : F \in \mathfrak{F}\}$$

is a continuous function because if $x \in S$, then there exists an open set W with $x \in W$ and W has nonempty intersection with only finitely many sets of $F \in \mathfrak{F}$. Then for $y \in W$,

$$\sum \{g_F(y) : F \in \mathfrak{F}\} = \sum_{i=1}^n g_{F_i}(y).$$

Since \mathfrak{F} is a cover of S ,

$$\sum \{g_F(x) : F \in \mathfrak{F}\} \neq 0$$

for any $x \in S$. Hence ϕ_F is continuous. This also shows $\phi_F(x) = 0$ for all but finitely many $F \in \mathfrak{F}$. It is obvious that

$$\sum \{\phi_F(x) : F \in \mathfrak{F}\} = 1$$

from the definition.

Let $z \in S$. Then there is an open set W containing z such that W has nonempty intersection with only finitely many $F \in \mathfrak{F}$. Thus for $y, x \in W$,

$$\left| \phi_{F_j}(x) - \phi_{F_j}(y) \right| \leq \left| \frac{g_{F_j}(x) \sum_{i=1}^n g_{F_i}(y) - g_{F_j}(y) \sum_{i=1}^n g_{F_i}(x)}{\sum_{i=1}^n g_{F_i}(x) \sum_{i=1}^n g_{F_i}(y)} \right|$$

If F is not one of these F_i , then $g_F(x) = \phi_F(x) = \phi_F(y) = g_F(y) = 0$. Thus there is nothing to show for these. It suffices to consider the ones above. Restricting W if necessary, we can assume that for $x \in W$,

$$\sum_F g_F(x) = \sum_{i=1}^n g_{F_i}(x) > \delta > 0, \quad g_{F_j}(x) < \Delta < \infty, \quad j \leq n$$

Then, simplifying the above, and letting $x, y \in W$, for each $j \leq n$,

$$\begin{aligned} \left| \phi_{F_j}(x) - \phi_{F_j}(y) \right| &\leq \frac{1}{\delta^2} \left| \begin{array}{l} g_{F_j}(x) \sum_F g_F(y) - g_{F_j}(y) \sum_F g_F(y) \\ + g_{F_j}(y) \sum_F g_F(y) - g_{F_j}(y) \sum_F g_F(x) \end{array} \right| \\ &\leq \frac{1}{\delta^2} \Delta |g_{F_j}(x) - g_{F_j}(y)| + \frac{1}{\delta^2} \Delta \sum_{i=1}^n |g_{F_i}(y) - g_{F_i}(x)| \\ &\leq \frac{\Delta}{\delta^2} d(x, y) + \frac{\Delta}{\delta^2} n d(x, y) = (n+1) \frac{\Delta}{\delta^2} d(x, y) \end{aligned}$$

Thus on this set W containing z , all ϕ_F are Lipschitz continuous with Lipschitz constant $(n+1) \frac{\Delta}{\delta^2}$. ■

The functions described above are called a partition of unity subordinate to the open cover \mathfrak{S} . A useful observation is contained in the following corollary.

Corollary 25.1.2 *Let S be a metric space and let \mathfrak{S} be any open cover of S . Then there exists a set \mathfrak{F} , an open refinement of \mathfrak{S} , and functions $\{\phi_F : F \in \mathfrak{F}\}$ such that*

$$\phi_F : S \rightarrow [0, 1]$$

ϕ_F is continuous

$\phi_F(x)$ equals 0 for all but finitely many $F \in \mathfrak{F}$

$$\sum \{\phi_F(x) : F \in \mathfrak{F}\} = 1 \text{ for all } x \in S.$$

Each ϕ_F is Lipschitz continuous. If $U \in \mathfrak{S}$ and H is a closed subset of U , the partition of unity can be chosen such that each $\phi_F = 0$ on H except for one which equals 1 on H .

Proof: Just change your open cover to consist of U and $V \setminus H$ for each $V \in \mathfrak{S}$. Then every function but one equals 0 on H and so exactly one of them equals 1 on H . ■

25.2 An Extension Theorem, Retracts

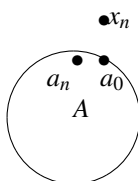
Recall the Tietze extension theorem which involved extending a real valued function defined on a closed set to a real valued function defined on the whole space. There is a big generalization in which the continuous function has values in a normed linear space. As with the Tietze extension theorem, the closed set is in a metric space.

Lemma 25.2.1 *Let A be a closed set in a metric space and let $x_n \notin A, x_n \rightarrow a_0 \in A$ and $a_n \in A$ such that $d(a_n, x_n) \leq 6 \text{dist}(x_n, A)$. Then $a_n \rightarrow a_0$.*

Proof: By assumption,

$$\begin{aligned} d(a_n, a_0) &\leq d(a_n, x_n) + d(x_n, a_0) < 6 \text{dist}(x_n, A) + d(x_n, a_0) \\ &\leq 6d(x_n, a_0) + d(x_n, a_0) = 7d(x_n, a_0) \end{aligned}$$

and this converges to 0. ■



In the proof of the following theorem, you get a covering of A^C with open balls B such that for each of these balls, there exists $a \in A$ such that for all $x \in B$, $\|x - a\|$ is no more than six times the distance of x to A . The 6 is not important. Any other constant with this property would work. Then you use Stone's theorem.

Recall a Banach space is a normed vector space which is also a complete metric space where the metric comes from the norm.

$$d(x, y) = \|x - y\|$$

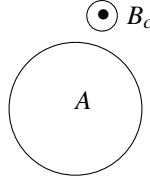
Thus you can add things in a Banach space.

Definition 25.2.2 A Banach space is a complete normed linear space. If you have a subset B of a Banach space, then $\text{conv}(B)$ denotes the smallest closed convex set which contains B . It can be obtained by taking the intersection of all closed convex sets containing B . Recall that a set C is convex if whenever $x, y \in C$, then so is $\lambda x + (1 - \lambda)y$ for all $\lambda \in [0, 1]$. Note how this makes sense in a vector space but maybe not in a general metric space.

In the following theorem, we have in mind both X and Y are Banach spaces, but this is not needed in the proof. All that is needed is that X is a metric space and Y a normed linear space or possibly something more general in which it makes sense to do addition and scalar multiplication.

Theorem 25.2.3 Let A be a closed subset of a metric space X and let $F : A \rightarrow Y$, Y a normed linear space. Then there exists an extension of F denoted as \hat{F} such that \hat{F} is defined on all of X and agrees with F on A . It has values in $\text{conv}(F(A))$, the convex hull of $F(A)$.

Proof: For each $c \notin A$, let B_c be a ball contained in A^C centered at c where distance of c to A is at least $\text{diam}(B_c)$.



So for $x \in B_c$ what about $\text{dist}(x, A)$? How does it compare with $\text{dist}(c, A)$?

$$\begin{aligned} \text{dist}(c, A) &\leq d(c, x) + \text{dist}(x, A) \leq \frac{1}{2} \text{diam}(B_c) + \text{dist}(x, A) \\ &\leq \frac{1}{2} \text{dist}(c, A) + \text{dist}(x, A) \end{aligned}$$

so $\text{dist}(c, A) \leq 2 \text{dist}(x, A)$. Now the following is also valid. Letting $x \in B_c$ be arbitrary, it follows from the assumption on the diameter that there exists $a_0 \in A$ such that $d(c, a_0) < 2 \text{dist}(c, A)$. Then

$$\begin{aligned} d(x, a_0) &\leq \sup_{y \in B_c} d(y, a_0) \leq \sup_{y \in B_c} (d(y, c) + d(c, a_0)) \leq \frac{\text{diam}(B_c)}{2} + 2 \text{dist}(c, A) \\ &\leq \frac{\text{dist}(c, A)}{2} + 2 \text{dist}(c, A) < 3 \text{dist}(c, A) \end{aligned} \quad (25.1)$$

It follows from 25.1, $d(x, a_0) \leq 3 \text{dist}(c, A) \leq 6 \text{dist}(x, A)$. Thus for any $x \in B_c$, there is an $a_0 \in A$ such that $d(x, a_0)$ is bounded by a fixed multiple of the distance from x to A .

By Stone's theorem, there is a locally finite open refinement \mathcal{R} . These are open sets each of which is contained in one of the balls just mentioned such that each of these balls is the union of sets of \mathcal{R} . Thus \mathcal{R} is a locally finite cover of A^C . Since $x \in A^C$ is in one of those balls, it was just shown that there exists $a_R \in A$ such that for all $x \in R \in \mathcal{R}$ we have $d(x, a_R) \leq 6 \text{dist}(x, A)$. Of course there may be more than one because R might be contained in more than one of those special balls. One a_R is chosen for each $R \in \mathcal{R}$.

Now let $\phi_R(x) \equiv \text{dist}(x, R^C)$. Then let

$$\hat{F}(x) \equiv \begin{cases} F(x) & \text{for } x \in A \\ \sum_{R \in \mathcal{R}} F(a_R) \frac{\phi_R(x)}{\sum_{\hat{R} \in \mathcal{R}} \phi_{\hat{R}}(x)} & \text{for } x \notin A \end{cases}$$

The sum in the bottom is always finite because the covering is locally finite. Also, this sum is never 0 because \mathcal{R} is a covering. Also \hat{F} has values in $\text{conv}(F(K))$. It only remains to verify that \hat{F} is continuous. It is clearly so on the interior of A thanks to continuity of F . It is also clearly continuous on A^C because the functions ϕ_R are continuous. So it suffices to consider $x_n \rightarrow a \in \partial A \subseteq A$ where $x_n \notin A$ and see whether $F(a) = \lim_{n \rightarrow \infty} \hat{F}(x_n)$.

Suppose this does not happen. Then there is a sequence converging to some $a \in \partial A$ and $\varepsilon > 0$ such that

$$\varepsilon \leq \|\hat{F}(a) - \hat{F}(x_n)\| \quad \text{all } n$$

For $x_n \in R$, it was shown above that $d(x_n, a_{R_n}) \leq 6 \text{dist}(x_n, A)$. By the above Lemma 25.2.1, it follows that $a_{R_n} \rightarrow a$ and so $F(a_{R_n}) \rightarrow F(a)$.

$$\varepsilon \leq \|\hat{F}(a) - \hat{F}(x_n)\| \leq \sum_{R \in \mathcal{R}} \|F(a_{R_n}) - F(a)\| \frac{\phi_R(x_{Rn})}{\sum_{\hat{R} \in \mathcal{R}} \phi_{\hat{R}}(x_{Rn})}$$

By local finiteness of the cover, each x_n involves only finitely many R . Thus, in this limit process, there are countably many R involved $\{R_j\}_{j=1}^\infty$. Thus one can apply Fatou's lemma.

$$\begin{aligned} \varepsilon &\leq \liminf_{n \rightarrow \infty} \|\hat{F}(a) - \hat{F}(x_n)\| \\ &\leq \sum_{j=1}^\infty \liminf_{n \rightarrow \infty} \|F(a_{R_{jn}}) - F(a)\| \frac{\phi_{R_j}(x_{R_{jn}})}{\sum_{j=1}^\infty \phi_{R_j}(x_{R_{jn}})} \\ &\leq \sum_{j=1}^\infty \liminf_{n \rightarrow \infty} \|F(a_{R_{jn}}) - F(a)\| = 0 \blacksquare \end{aligned}$$

The last step is needed because you lose local finiteness as you approach ∂A . Note that the only thing needed was that X is a metric space. The addition takes place in Y so it needs to be a vector space. Did it need to be complete? No, this was not used. Nor was completeness of X used. The main interest here is in Banach spaces, but the result is more general than that.

It also appears that \hat{F} is locally Lipschitz on A^C .

Definition 25.2.4 Let S be a subset of X , a Banach space. Then it is a retract if there exists a continuous function $R : X \rightarrow S$ such that $Rs = s$ for all $s \in S$. This R is a retraction. More generally, $S \subseteq T$ is called a retract of T if there is a continuous $R : T \rightarrow S$ such that $Rs = s$ for all $s \in S$.

Theorem 25.2.5 Let K be closed and convex subset of X a Banach space. Then K is a retract.

Proof: By Theorem 25.2.3, there is a continuous function \hat{I} extending I to all of X . Then also \hat{I} has values in $\text{conv}(IK) = \text{conv}(K) = K$. Hence \hat{I} is a continuous function which does what is needed. It maps everything into K and keeps the points of K unchanged. ■

Sometimes people call the set a retraction also or the function which does the job a retraction. This seems like strange thing to call it because a retraction is the act of repudiating something you said earlier. Nevertheless, I will call it that. Note that if S is a retract of the whole metric space X , then it must be a retract of every set which contains S .

Part IV

**Stochastic Processes and
Probability**

Chapter 26

Independence

Caution: This material on probability and stochastic processes may be half baked in places. This is not to say that nothing else is half baked. However, the probability is higher here. Probability is not my main research area so my qualifications for even writing this are not all they could be. However, I like probability and think it fits in well with what is presented earlier and hope this might be useful for someone like me. This book is not a research monograph written for experts, and I am no expert. I have included mainly those items which I have found most interesting. However, there is an awful lot in this subject, far more than I can include. For more topics see the references.

This material was written down earlier in the Topics in Analysis book but in a haphazard manner as I encountered it rather than in the most logical manner. I am trying to present it here in a more coherent form.

26.1 Random Variables and Independence

Recall Lemma 20.2.3 on Page 526 which is stated here for convenience.

Lemma 26.1.1 *Let M be a metric space with the closed balls compact and suppose λ is a measure defined on the Borel sets of M which is finite on compact sets. Then there exists a unique Radon measure, $\bar{\lambda}$ which equals λ on the Borel sets. In particular λ must be both inner and outer regular on all Borel sets.*

Also important is the following fundamental result which is called the Borel Cantelli lemma. It is Lemma 9.2.5 on Page 243.

Lemma 26.1.2 *Let $(\Omega, \mathcal{F}, \lambda)$ be a measure space and let $\{A_i\}$ be a sequence of measurable sets satisfying $\sum_{i=1}^{\infty} \lambda(A_i) < \infty$. Then letting S denote the set of $\omega \in \Omega$ which are in infinitely many A_i , it follows S is a measurable set and $\lambda(S) = 0$.*

Here is another nice observation.

Proposition 26.1.3 *Suppose E_i is a separable Banach space. Then if B_i is a Borel set of E_i , it follows $\prod_{i=1}^n B_i$ is a Borel set in $\prod_{i=1}^n E_i$.*

Proof: An easy way to do this is to consider the projection maps $\pi_i x \equiv x_i$. Then these projection maps are continuous. Hence for U open, $\pi_i^{-1}(U) \equiv \prod_{j=1}^n A_j$, $A_j = E_j$ if $j \neq i$ and $A_i = U$. Thus $\pi_i^{-1}(\text{open})$ equals an open set. Let $\mathcal{S} \equiv \{V \subseteq \mathbb{R} : \pi_i^{-1}(V) \text{ is Borel}\}$. Then \mathcal{S} contains all the open sets and is clearly a σ algebra. Therefore, \mathcal{S} contains the Borel sets. Let B_i be a Borel set in E_i . Then $\prod_{i=1}^n B_i = \cap_{i=1}^n \pi_i^{-1}(B_i)$, a finite intersection of Borel sets. ■

Definition 26.1.4 *A probability space is a measure space, (Ω, \mathcal{F}, P) where P is a measure satisfying $P(\Omega) = 1$. A random vector (variable) is a measurable function, $\mathbf{X} : \Omega \rightarrow Z$ where Z is some topological space. It is often the case that Z will equal \mathbb{R}^p . Assume Z is a separable Banach space. Define the following σ algebra.*

$$\sigma(\mathbf{X}) \equiv \{\mathbf{X}^{-1}(E) : E \text{ is Borel in } Z\}$$

Thus $\sigma(\mathbf{X}) \subseteq \mathcal{F}$. For E a Borel set in Z define $\lambda_{\mathbf{X}}(E) \equiv P(\mathbf{X}^{-1}(E))$. This is called the distribution of the random variable \mathbf{X} . If $\int_{\Omega} |\mathbf{X}(\omega)| dP < \infty$ then define $E(\mathbf{X}) \equiv \int_{\Omega} \mathbf{X} dP$ where the integral is defined as the Bochner integral.

Recall the following fundamental result which was proved earlier but which I will give a short proof of now. Recall Definition 24.1.1 about strongly measurable being the limit of simple functions.

Proposition 26.1.5 *Let $(\Omega, \mathcal{S}, \mu)$ be a measure space and let $X : \Omega \rightarrow Z$ where Z is a separable Banach space. Then X is strongly measurable if and only if $X^{-1}(U) \in \mathcal{S}$ for all U open in Z .*

Proof: To begin with, let $D(a, r)$ be the closure of the open ball $B(a, r)$. By Lemma 24.1.7, there exists $\{f_i\} \subseteq B'$, the unit ball in Z' such that $\|z\|_Z = \sup_i \{|f_i(z)|\}$. Then

$$\begin{aligned} D(a, r) &= \{z : \|a - z\| \leq r\} = \cap_i \{z : |f_i(z) - f_i(a)| \leq r\} \\ &= \cap_i f_i^{-1} \left(\overline{B(f_i(a), r)} \right) \end{aligned}$$

Thus $X^{-1}(D(a, r)) = \cap_i X^{-1} \left(f_i^{-1} \left(\overline{B(f_i(a), r)} \right) \right) = \cap_i (f_i \circ X)^{-1} \left(\overline{B(f_i(a), r)} \right)$. If X is strongly measurable, then it is weakly measurable and so each $f_i \circ X$ is a real (complex) valued measurable function. Hence the expression on the right in the above is measurable. Now if U is any open set in Z , then it is the countable union of such closed disks $U = \cup_i D_i$. Therefore, $X^{-1}(U) = \cap_i X^{-1}(D_i) \in \mathcal{S}$. It follows that strongly measurable implies inverse images of open sets are in \mathcal{S} .

Conversely, suppose $X^{-1}(U) \in \mathcal{S}$ for every open U . Then for $f \in Z'$, $f \circ X$ is real valued and measurable. Therefore, X is weakly measurable. By the Pettis theorem, Theorem 24.1.8, it follows that $f \circ X$ is strongly measurable. ■

Proposition 26.1.6 *If $X : \Omega \rightarrow Z$ is measurable, then $\sigma(X)$ equals the smallest σ algebra such that X is measurable with respect to it. Also if X_i are random variables having values in separable Banach spaces Z_i , then $\sigma(X) = \sigma(X_1, \dots, X_n)$ where X is the vector mapping Ω to $\prod_{i=1}^n Z_i$ and $\sigma(X_1, \dots, X_n)$ is the smallest σ algebra such that each X_i is measurable with respect to it.*

Proof: Let \mathcal{G} denote the smallest σ algebra such that X is measurable with respect to this σ algebra. By definition $X^{-1}(\text{open}) \in \mathcal{G}$. Furthermore, the set of all E such that $X^{-1}(E) \in \mathcal{G}$ is a σ algebra. Hence it includes all the Borel sets. Hence $X^{-1}(\text{Borel}) \in \mathcal{G}$ and so $\mathcal{G} \supseteq \sigma(X)$. However, $\sigma(X)$ defined above is a σ algebra such that X is measurable with respect to $\sigma(X)$. Therefore, $\mathcal{G} = \sigma(X)$.

Letting B_i be a Borel set in Z_i , $\prod_{i=1}^n B_i$ is a Borel set by Proposition 26.1.3 and so $X^{-1}(\prod_{i=1}^n B_i) = \cap_{i=1}^n X_i^{-1}(B_i) \in \sigma(X_1, \dots, X_n)$. If \mathcal{G} denotes the Borel sets $F \subseteq \prod_{i=1}^n Z_i$ such that $X^{-1}(F) \in \sigma(X_1, \dots, X_n)$, then \mathcal{G} is clearly a σ algebra which contains the open sets. Hence $\mathcal{G} = \mathcal{B}$ the Borel sets of $\prod_{i=1}^n Z_i$. This shows that $\sigma(X) \subseteq \sigma(X_1, \dots, X_n)$. Next we observe that $\sigma(X)$ is a σ algebra with the property that each X_i is measurable with respect to $\sigma(X)$. This follows from $X_i^{-1}(B_i) = X^{-1}(\prod_{j=1}^n A_j) \in \sigma(X)$, where each $A_j = Z_j$ except for $A_i = B_i$. Since $\sigma(X_1, \dots, X_n)$ is defined as the smallest such σ algebra, it follows that $\sigma(X) \supseteq \sigma(X_1, \dots, X_n)$. ■

For random variables having values in a separable Banach space or even more generally for a separable metric space, much can be said about regularity of λ_X .

Definition 26.1.7 A measure, μ defined on $\mathcal{B}(E)$ for E a separable metric space will be called inner regular if for all $F \in \mathcal{B}(E)$,

$$\mu(F) = \sup\{\mu(K) : K \subseteq F \text{ and } K \text{ is closed}\}$$

A measure, μ defined on $\mathcal{B}(E)$ will be called outer regular if for all $F \in \mathcal{B}(E)$,

$$\mu(F) = \inf\{\mu(V) : V \supseteq F \text{ and } V \text{ is open}\}$$

When a measure is both inner and outer regular, it is called regular.

Note that if the metric space is \mathbb{R}^p then λ_X can be considered a Radon measure because you can use it to obtain a positive linear functional and then use the Riesz representation theorem for these.

For probability measures, the above definition of regularity tends to come free. Note it is a little weaker than the usual definition of regularity because K is only assumed to be closed, not compact. This is stated for convenience. It is Lemma 9.8.4 on Page 253.

Lemma 26.1.8 Let μ be a finite measure defined on $\mathcal{B}(E)$ where E is a metric space. Then μ is regular.

One can say more if the metric space is complete and separable. In fact in this case the above definition of inner regularity can be shown to imply the usual one where the closed sets are replaced with compact sets. It is Lemma 9.8.5 on Page 255.

Lemma 26.1.9 Let μ be a finite measure on a σ algebra containing $\mathcal{B}(X)$, the Borel sets of X , a separable complete metric space. (Polish space) Then if C is a closed set,

$$\mu(C) = \sup\{\mu(K) : K \subseteq C \text{ and } K \text{ is compact}\}$$

It follows that for a finite measure on $\mathcal{B}(X)$ where X is a Polish space, μ is inner regular in the sense that for all $F \in \mathcal{B}(X)$,

$$\mu(F) = \sup\{\mu(K) : K \subseteq F \text{ and } K \text{ is compact}\}$$

Definition 26.1.10 A measurable function $X : (\Omega, \mathcal{F}, \mu) \rightarrow Z$ a topological space is called a random variable when $\mu(\Omega) = 1$. For such a random variable, one can define a distribution measure λ_X on the Borel sets of Z as follows: $\lambda_X(G) \equiv \mu(X^{-1}(G))$. This is a well defined measure on the Borel sets of Z because it makes sense for every G open and $\mathcal{G} \equiv \{G \subseteq Z : X^{-1}(G) \in \mathcal{F}\}$ is a σ algebra which contains the open sets, hence the Borel sets. Such a measurable function is also called a random vector.

Corollary 26.1.11 Let X be a random variable (random vector) with values in a complete metric space, Z . Then λ_X is an inner and outer regular measure defined on $\mathcal{B}(Z)$.

Proposition 26.1.12 For X a random vector defined above, X having values in a complete separable metric space Z , then λ_X is inner and outer regular and Borel.

$$(\Omega, P) \xrightarrow{X} (Z, \lambda_X) \xrightarrow{h} E$$

If h is Borel measurable and $h \in L^1(Z, \lambda_X; E)$ for E a Banach space, then

$$\int_{\Omega} h(X(\omega)) dP = \int_Z h(x) d\lambda_X. \quad (26.1)$$

In the case where $Z = E$, a separable Banach space, if \mathbf{X} is measurable then $\mathbf{X} \in L^1(\Omega; E)$ if and only if the identity map on E is in $L^1(E; \lambda_{\mathbf{X}})$ and

$$\int_{\Omega} \mathbf{X}(\omega) dP = \int_E x d\lambda_{\mathbf{X}}(x) \quad (26.2)$$

Proof: The regularity claims are established above. It remains to verify 26.1.

Since $h \in L^1(Z, E)$, it follows there exists a sequence of simple functions $\{h_n\}$ such that

$$h_n(x) \rightarrow h(x), \quad \int_Z \|h_m - h_n\| d\lambda_{\mathbf{X}} \rightarrow 0 \text{ as } m, n \rightarrow \infty.$$

The first convergence above implies

$$h_n \circ \mathbf{X} \rightarrow h \circ \mathbf{X} \text{ pointwise on } \Omega \quad (26.3)$$

Then letting $h_n(x) = \sum_{k=1}^m x_k \mathcal{X}_{E_k}(x)$, where the E_k are disjoint and Borel, it follows easily that $h_n \circ \mathbf{X}$ is also a simple function of the form $h_n \circ \mathbf{X}(\omega) = \sum_{k=1}^m x_k \mathcal{X}_{\mathbf{X}^{-1}(E_k)}(\omega)$ and by assumption $\mathbf{X}^{-1}(E_k) \in \mathcal{F}$. From the definition of the integral, it is easily seen

$$\int h_n \circ \mathbf{X} dP = \int h_n d\lambda_{\mathbf{X}}, \quad \int \|h_n\| \circ \mathbf{X} dP = \int \|h_n\| d\lambda_{\mathbf{X}}$$

Also, $h_n \circ \mathbf{X} - h_m \circ \mathbf{X}$ is a simple function and so

$$\int \|h_n \circ \mathbf{X} - h_m \circ \mathbf{X}\| dP = \int \|h_n - h_m\| d\lambda_{\mathbf{X}} \quad (26.4)$$

It follows from the definition of the Bochner integral and 26.3, and 26.4 that $h \circ \mathbf{X}$ is in $L^1(\Omega; E)$ and

$$\int h \circ \mathbf{X} dP = \lim_{n \rightarrow \infty} \int h_n \circ \mathbf{X} dP = \lim_{n \rightarrow \infty} \int h_n d\lambda_{\mathbf{X}} = \int h d\lambda_{\mathbf{X}}.$$

Finally consider the case that $E = Z$ for E a separable Banach space, and suppose $\mathbf{X} \in L^1(\Omega; E)$. Then letting h be the identity map on E , it follows h is obviously separably valued and $h^{-1}(U) \in \mathcal{B}(E)$ for all U open and so h is measurable. Why is it in $L^1(E; E)$?

$$\begin{aligned} \int_E \|h(x)\| d\lambda_{\mathbf{X}} &= \int_0^{\infty} \lambda_{\mathbf{X}}(\{\|h\| > t\}) dt \equiv \int_0^{\infty} P(\mathbf{X} \in \{\|x\| > t\}) dt \\ &\equiv \int_0^{\infty} P(\|\mathbf{X}\| > t) dt = \int_{\Omega} \|\mathbf{X}\| dP < \infty \end{aligned}$$

Thus the identity map on E is in $L^1(E; \lambda_{\mathbf{X}})$. Next let the identity map h be in $L^1(E; \lambda_{\mathbf{X}})$. Then $\mathbf{X}(\omega) = h \circ \mathbf{X}(\omega)$ and so from the first part, $\mathbf{X} \in L^1(\Omega; E)$ and from 26.1, 26.2 follows. ■

26.2 Convergence in Probability

Definition 26.2.1 $\{f_n\}$ is said to be Cauchy in probability if for each $\varepsilon > 0$,

$$\lim_{n, m \rightarrow \infty} P(\|f_n - f_m\| > \varepsilon) = 0$$

This means: for each $\delta > 0$ there exists k_{δ} such that if $m, n \geq k_{\delta}$, then $P(\|f_n - f_m\| > \varepsilon) < \delta$.

Proposition 26.2.2 $\{f_n\}$ is Cauchy in probability, these being functions having values in a Banach space then there exists a set of measure zero N and a subsequence $\{f_{n_k}\}$ such that for $\omega \notin N$, $\lim_{k \rightarrow \infty} f_{n_k}(\omega)$ converges.

Proof: From the above definition, there exists n_1 such that if $m, n \geq n_1$, then

$$P(\|f_n - f_m\| > 2^{-1}) < 2^{-1}$$

From the definition, there exists $n_2 > n_1$ such that $P(\|f_n - f_m\| > 2^{-2}) < 2^{-2}$ whenever $n, m \geq n_1$ and so forth. Thus

$$P(\|f_{n_{k+1}} - f_{n_k}\| > 2^{-k}) < 2^{-k}$$

Letting $A_k \equiv [\|f_{n_{k+1}} - f_{n_k}\| > 2^{-k}]$, it follows from the Borell Cantelli lemma that there is a set of measure zero, namely $N \equiv \bigcap_{n=1}^{\infty} \bigcup_{k \geq n} A_k$ such that $P(N^C) = 1 = P(\bigcup_{n=1}^{\infty} \bigcap_{k \geq n} A_k^C)$. To say $\omega \in N^C$ is the same as saying that there exists n such that ω is in A_k^C for all $k \geq n$. In other words, eventually $\|f_{n_{k+1}} - f_{n_k}\| \leq 2^{-k}$. Now it follows that

$$\|f_{n+p}(\omega) - f_{n_k}(\omega)\| \leq \sum_{j=k}^{\infty} \|f_{n_{j+1}}(\omega) - f_{n_j}(\omega)\| < 2^{-(k-1)}$$

if k is large enough. Hence $\{f_{n_k}(\omega)\}_k$ is a Cauchy sequence for each $\omega \notin N$ and since E is complete, this sequence converges. ■

26.3 Kolmogorov Extension Theorem

Let M_t be a complete separable metric space. This is called a Polish space. I will denote a totally ordered index set, (Like \mathbb{R}) and the interest will be in building a measure on the product space, $\prod_{t \in I} M_t$. If you like less generality, just think of $M_t = \mathbb{R}^{k_t}$ or even $M_t = \mathbb{R}$. By the well ordering principle, you can always put an order on any index set so this order is no restriction, but we do not insist on a well order and in fact, index sets of great interest are \mathbb{R} or $[0, \infty)$. Also for X a topological space, $\mathcal{B}(X)$ will denote the Borel sets.

Notation 26.3.1 The symbol J will denote a finite subset of I , $J = (t_1, \dots, t_n)$, the t_i taken in order. E_J will denote a set which has a set E_t of $\mathcal{B}(M_t)$ in the t^{th} position for $t \in J$ and for $t \notin J$, the set in the t^{th} position will be M_t . K_J will denote a set which has a compact set in the t^{th} position for $t \in J$ and for $t \notin J$, the set in the t^{th} position will be M_t . Also denote by \mathcal{R}_J the sets E_J and \mathcal{R} the union of all such \mathcal{R}_J . Let \mathcal{E}_J denote finite disjoint unions of sets of \mathcal{R}_J and let \mathcal{E} denote finite disjoint unions of sets of \mathcal{R} . Thus if F is a set of \mathcal{E} , there exists J such that F is a finite disjoint union of sets of \mathcal{R}_J . For $F \in \Omega$, denote by $\pi_J(F)$ the set $\prod_{t \in J} F_t$ where $F = \prod_{t \in I} F_t$.

With this preparation, here is the Kolmogorov extension theorem. It is Theorem 20.3.3 proved earlier. In the statement and proof of the theorem, F_i, G_i , and E_i will denote Borel sets. Any list of indices from I will always be assumed to be taken in order. Thus, if $J \subseteq I$ and $J = (t_1, \dots, t_n)$, it will always be assumed $t_1 < t_2 < \dots < t_n$.

Theorem 26.3.2 For each finite set $J = (t_1, \dots, t_n) \subseteq I$, suppose there exists a Borel probability measure, $\nu_J = \nu_{t_1 \dots t_n}$ defined on the Borel sets of $\prod_{t \in J} M_t$ such that the following consistency condition holds. If $(t_1, \dots, t_n) \subseteq (s_1, \dots, s_p)$, then

$$\nu_{t_1 \dots t_n}(F_{t_1} \times \dots \times F_{t_n}) = \nu_{s_1 \dots s_p}(G_{s_1} \times \dots \times G_{s_p}) \quad (26.5)$$

where if $s_i = t_j$, then $G_{s_i} = F_{t_j}$ and if s_i is not equal to any of the indices t_k , then $G_{s_i} = M_{s_i}$. Then for \mathcal{E} defined in the above Notation, there exists a probability measure P and a σ algebra $\mathcal{F} = \sigma(\mathcal{E})$ such that $(\prod_{t \in I} M_t, P, \mathcal{F})$ is a probability space. Also there exist measurable functions, $X_s : \prod_{t \in I} M_t \rightarrow M_s$ defined for $s \in I$ as

$$X_s \mathbf{x} \equiv x_s$$

such that for each $(t_1 \cdots t_n) \subseteq I$,

$$\begin{aligned} \nu_{t_1 \cdots t_n}(F_{t_1} \times \cdots \times F_{t_n}) &= P([X_{t_1} \in F_{t_1}] \cap \cdots \cap [X_{t_n} \in F_{t_n}]) \\ &= P\left((X_{t_1}, \cdots, X_{t_n}) \in \prod_{j=1}^n F_{t_j}\right) = P\left(\prod_{t \in I} F_t\right) \end{aligned} \quad (26.6)$$

where $F_t = M_t$ for every $t \notin \{t_1 \cdots t_n\}$ and F_{t_i} is a Borel set. Also if f is a nonnegative function of finitely many variables, x_{t_1}, \cdots, x_{t_n} , measurable with respect to $\mathcal{B}\left(\prod_{j=1}^n M_{t_j}\right)$, then f is also measurable with respect to \mathcal{F} and

$$\int_{M_{t_1} \times \cdots \times M_{t_n}} f(x_{t_1}, \cdots, x_{t_n}) d\nu_{t_1 \cdots t_n} = \int_{\prod_{t \in I} M_t} f(x_{t_1}, \cdots, x_{t_n}) dP \quad (26.7)$$

26.4 Independent Events and σ Algebras

The concept of independence is probably the main idea which separates probability from analysis and causes some of us, myself included, to struggle to understand what is going on. I think that these ideas are the main difficulty some encounter when trying to understand probability, not the kind based on combinatorics but what is being presented here.

Definition 26.4.1 Let (Ω, \mathcal{F}, P) be a probability space. The sets in \mathcal{F} are called events. A set of events, $\{A_i\}_{i \in I}$ is called independent if whenever $\{A_{i_k}\}_{k=1}^m$ is a finite subset

$$P\left(\bigcap_{k=1}^m A_{i_k}\right) = \prod_{k=1}^m P(A_{i_k}).$$

Each of these events defines a rather simple σ algebra, $(A_i, A_i^C, \emptyset, \Omega)$ denoted by \mathcal{F}_i . Now the following lemma is interesting because it motivates a more general notion of independent σ algebras.

Lemma 26.4.2 Suppose $\{A_i\}_{i \in I}$ are independent events. Then for

$$B_i \in \mathcal{F}_i \equiv \{A_i, A_i^C, \emptyset, \Omega\}$$

for $i \in I$. Then for any $m \in \mathbb{N}$, $P\left(\bigcap_{k=1}^m B_{i_k}\right) = \prod_{k=1}^m P(B_{i_k})$.

Proof: The proof is by induction on the number l of the B_{i_k} which are not equal to A_{i_k} . First suppose $l = 0$. Then the above assertion is true by assumption since $\{A_i\}_{i \in I}$ is independent. Suppose it is so for some l and there are $l + 1$ sets not equal to A_{i_k} . If any equals \emptyset there is nothing to show. Both sides equal 0. If any equals Ω , there is also nothing to show. You can ignore that set in both sides and then you have by induction the two sides

are equal because you have no more than l sets different than A_{i_k} . The only remaining case is where some $B_{i_k} = A_{i_k}^C$. Say $B_{i_{m+1}} = A_{i_{m+1}}^C$ for simplicity.

$$P\left(\bigcap_{k=1}^{m+1} B_{i_k}\right) = P\left(A_{i_{m+1}}^C \cap \bigcap_{k=1}^m B_{i_k}\right) = P\left(\bigcap_{k=1}^m B_{i_k}\right) - P\left(A_{i_{m+1}} \cap \bigcap_{k=1}^m B_{i_k}\right)$$

Then by induction,

$$\begin{aligned} &= \prod_{k=1}^m P(B_{i_k}) - P(A_{i_{m+1}}) \prod_{k=1}^m P(B_{i_k}) = \prod_{k=1}^m P(B_{i_k}) (1 - P(A_{i_{m+1}})) \\ &= P(A_{i_{m+1}}^C) \prod_{k=1}^m P(B_{i_k}) = \prod_{k=1}^{m+1} P(B_{i_k}) \end{aligned}$$

thus proving it for $l+1$. ■

This motivates a more general notion of independence in terms of σ algebras.

Definition 26.4.3 If $\{\mathcal{F}_i\}_{i \in I}$ is any set of σ algebras contained in \mathcal{F} , they are said to be independent if whenever $A_{i_k} \in \mathcal{F}_{i_k}$ for $k = 1, 2, \dots, m$, then

$$P\left(\bigcap_{k=1}^m A_{i_k}\right) = \prod_{k=1}^m P(A_{i_k}).$$

A set of random variables $\{X_i\}_{i \in I}$ is independent if the σ algebras $\{\sigma(X_i)\}_{i \in I}$ are independent σ algebras. Here $\sigma(X)$ denotes the smallest σ algebra such that X is measurable. Thus $\sigma(X) = \{X^{-1}(U) : U \text{ is a Borel set}\}$. More generally, $\sigma(X_i : i \in I)$ is the smallest σ algebra such that each X_i is measurable.

Note that by Lemma 26.4.2 you can consider independent events in terms of independent σ algebras. That is, a set of independent events can always be considered as events taken from a set of independent σ algebras. This is a more general notion because here the σ algebras might have infinitely many sets in them.

Lemma 26.4.4 Suppose the set of random variables, $\{X_i\}_{i \in I}$ is independent. Also suppose $I_1 \subseteq I$ and $j \notin I_1$. Then the σ algebras $\sigma(X_i : i \in I_1)$, $\sigma(X_j)$ are independent σ algebras.

Proof: Let $B \in \sigma(X_j)$. I want to show that for any $A \in \sigma(X_i : i \in I_1)$, it follows that $P(A \cap B) = P(A)P(B)$. Let \mathcal{K} consist of finite intersections of sets of the form $X_k^{-1}(B_k)$ where B_k is a Borel set and $k \in I_1$. Thus \mathcal{K} is a π system and $\sigma(\mathcal{K}) = \sigma(X_i : i \in I_1)$. This is because it follows from the definition that $\sigma(\mathcal{K}) \supseteq \sigma(X_i : i \in I_1)$ because $\sigma(\mathcal{K})$ contains all $X_i^{-1}(B)$ for B Borel. For the other inclusion, the right side consists of all sets $X_i^{-1}(B)$ where B is a Borel set and so the right side, being a σ algebra also must contain all finite intersections of these sets which means $\sigma(X_i : i \in I_1)$ must contain \mathcal{K} and so $\sigma(X_i : i \in I_1) \supseteq \sigma(\mathcal{K})$.

Now if you have one of these sets of the form $A = \bigcap_{k=1}^m X_k^{-1}(B_k)$ where without loss of generality, it can be assumed the k are distinct since $X_k^{-1}(B_k) \cap X_k^{-1}(B'_k) = X_k^{-1}(B_k \cap B'_k)$, then

$$\begin{aligned} P(A \cap B) &= P\left(\bigcap_{k=1}^m X_k^{-1}(B_k) \cap B\right) = P(B) \prod_{k=1}^m P(X_k^{-1}(B_k)) \\ &= P(B) P\left(\bigcap_{k=1}^m X_k^{-1}(B_k)\right). \end{aligned}$$

Thus \mathcal{K} is contained in

$$\mathcal{G} \equiv \{A \in \sigma(\mathbf{X}_i : i \in I_1) : P(A \cap B) = P(A)P(B)\}.$$

Now \mathcal{G} is closed with respect to complements and countable disjoint unions. Here is why: If each $A_i \in \mathcal{G}$ and the A_i are disjoint,

$$\begin{aligned} P((\cup_{i=1}^{\infty} A_i) \cap B) &= P(\cup_{i=1}^{\infty} (A_i \cap B)) \\ &= \sum_i P(A_i \cap B) = \sum_i P(A_i)P(B) \\ &= P(B) \sum_i P(A_i) = P(B)P(\cup_{i=1}^{\infty} A_i) \end{aligned}$$

If $A \in \mathcal{G}$, $P(A^C \cap B) + P(A \cap B) = P(B)$ and so

$$\begin{aligned} P(A^C \cap B) &= P(B) - P(A \cap B) = P(B) - P(A)P(B) \\ &= P(B)(1 - P(A)) = P(B)P(A^C). \end{aligned}$$

Therefore, from the lemma on π systems, Lemma 9.3.2 on Page 243, it follows

$$\sigma(\mathbf{X}_i : i \in I_1) \supseteq \mathcal{G} \supseteq \sigma(\mathcal{K}) = \sigma(\mathbf{X}_i : i \in I_1). \blacksquare$$

Definition 26.4.5 When X is a random variable with values in a Banach space Z , the notation $E(X)$ means $\int_Z X(\omega) dP$ where the latter is the Bochner integral. I will use this notation whenever convenient. $E(X)$ is called the expectation of X or the expected value of X . I will sometimes also use E as the name of a Banach space, but it should be clear from the context which is meant.

Recall Lemma 10.16.1

Lemma 26.4.6 Let f, g be nonnegative measurable nonnegative functions on a measure space (Ω, μ) . Then $\int f g d\mu = \int_0^{\infty} \int_{[g>t]} f d\mu dt = \int_0^{\infty} \int_0^{\infty} \mu([f>s] \cap [g>t]) ds dt$.

Corollary 26.4.7 If $\{f_i\}_{i=1}^m$ are nonnegative measurable functions, it follows from induction that

$$\int \prod_{i=1}^m f_i d\mu = \int_0^{\infty} \cdots \int_0^{\infty} \mu(\cap_{i=1}^m [f_i > t_i]) dt_1 \cdots dt_m$$

Proof: The case of $n = 2$ was just done. So suppose true for $n \geq 2$. Then from this case and induction,

$$\begin{aligned} \int \prod_{i=1}^{n+1} f_i d\mu &= \int_0^{\infty} \int_{[f_{n+1} > t_{n+1}]} \prod_{i=1}^n f_i d\mu dt_{n+1} = \int_0^{\infty} \int \prod_{i=1}^n \mathcal{X}_{[f_{n+1} > t_{n+1}]} f_i d\mu dt_{n+1} \\ &= \int_0^{\infty} \int_0^{\infty} \cdots \int_0^{\infty} \mu(\cap_{i=1}^n [\mathcal{X}_{[f_{n+1} > t_{n+1}]} f_i > t_i]) dt_1 \cdots dt_n dt_{n+1} \\ &= \int_0^{\infty} \cdots \int_0^{\infty} \mu(\cap_{i=1}^n [f_{n+1} > t_{n+1}] \cap [f_i > t_i]) dt_1 \cdots dt_{n+1} \\ &= \int_0^{\infty} \cdots \int_0^{\infty} \mu(\cap_{i=1}^{n+1} [f_i > t_i]) dt_1 \cdots dt_{n+1} \blacksquare \end{aligned}$$

Lemma 26.4.8 If $\{X_k\}_{k=1}^r$ are independent random variables having values in Z a separable metric space, and if g_k is a Borel measurable function, then $\{g_k(X_k)\}_{k=1}^r$ is also independent.

Proof: First consider the claim about $\{g_k(X_k)\}_{k=1}^r$. Letting O be an open set in Z ,

$$(g_k \circ X_k)^{-1}(O) = X_k^{-1}(g_k^{-1}(O)) = X_k^{-1}(\text{Borel set}) \in \sigma(X_k).$$

It follows $(g_k \circ X_k)^{-1}(E)$ is in $\sigma(X_k)$ whenever E is Borel because the sets whose inverse images are measurable includes the Borel sets. Thus $\sigma(g_k \circ X_k) \subseteq \sigma(X_k)$. ■

Theorem 26.4.9 Suppose $\{X_i\}_{i=1}^m$ are independent random variables with values in X a Banach space, then $\prod_{i=1}^m \|X_i\| \int_{\Omega} \prod_{i=1}^m \|X_i\| dP = \prod_{i=1}^m \int_{\Omega} \|X_i\| dP$.

Proof: The real valued random variables $\|X_i\|$ are respectively measurable in $\sigma(X_i)$ and so, from Corollary 26.4.7 and the independence of $\|X_i\|$,

$$\begin{aligned} \int \prod_{i=1}^m \|X_i\| d\mu &= \int_0^\infty \cdots \int_0^\infty \mu(\cap_{i=1}^m [\|X_i\|_i > t_i]) dt_1 \cdots dt_m \\ &= \int_0^\infty \cdots \int_0^\infty \prod_{i=1}^m \mu([\|X_i\|_i > t_i]) dt_1 \cdots dt_m \\ &= \prod_{i=1}^m \int_0^\infty \mu([\|X_i\|_i > t_i]) dt_i = \prod_{i=1}^m \int_{\Omega} \|X_i\| dP \quad \blacksquare \end{aligned}$$

Note that if each $X_i \in L^1$ and these are independent, then their product is also in L^1 .

Maybe this would be a good place to put a really interesting result known as the Doob Dynkin lemma. This amazing result is illustrated with the following diagram in which $\mathbf{X} = (X_1, \dots, X_m)$. By Proposition 26.1.6 $\sigma(\mathbf{X}) = \sigma(X_1, \dots, X_m)$, the expression on the right being the smallest σ algebra such that each X_i is measurable. The following diagram summarizes this result.

$$\begin{array}{ccc} (\Omega, \sigma(\mathbf{X})) & \xrightarrow{Y} & F \\ \searrow \mathbf{X} & \circ & \nearrow g \\ & (\prod_{i=1}^m E_i, \mathcal{B}(\prod_{i=1}^m E_i)) & \end{array}$$

You start with Y and can write it as the composition $g \circ \mathbf{X}$ provided Y is $\sigma(\mathbf{X})$ measurable.

Lemma 26.4.10 Let (Ω, \mathcal{F}) be a measure space and let $X_i : \Omega \rightarrow E_i$ where E_i is a separable Banach space. Suppose also that $Y : \Omega \rightarrow F$ where F is a separable Banach space. Then Y is $\sigma(X_1, \dots, X_m)$ measurable if and only if there exists a Borel measurable function $g : \prod_{i=1}^m E_i \rightarrow F$ such that $Y = g(X_1, \dots, X_m)$.

Proof: First suppose $Y(\omega) = f \mathcal{X}_W(\omega)$ where $f \in F$ and $W \in \sigma(X_1, \dots, X_m)$. Then by Proposition 26.1.6, W is of the form $(X_1, \dots, X_m)^{-1}(B) \equiv \mathbf{X}^{-1}(B)$ where B is Borel in $\prod_{i=1}^m E_i$. Therefore,

$$Y(\omega) = f \mathcal{X}_{\mathbf{X}^{-1}(B)}(\omega) = f \mathcal{X}_B(\mathbf{X}(\omega)).$$

Now suppose Y is measurable with respect to $\sigma(X_1, \dots, X_m)$. Then there exist simple functions

$$Y_n(\omega) = \sum_{k=1}^{m_n} f_k \mathcal{X}_{B_k}(X(\omega)) \equiv g_n(X(\omega))$$

where the B_k are Borel sets in $\prod_{i=1}^m E_i$, such that $Y_n(\omega) \rightarrow Y(\omega)$, each g_n being Borel. Thus g_n converges on $X(\Omega)$. Furthermore, the set on which g_n does converge is a Borel set equal to

$$\bigcap_{n=1}^{\infty} \bigcup_{m=1}^{\infty} \bigcap_{p,q \geq m} \left[\|g_p - g_q\| < \frac{1}{n} \right]$$

which contains $X(\Omega)$. Therefore, modifying g_n by multiplying it by the indicator function of this Borel set containing $X(\Omega)$, we can conclude that g_n converges to a Borel function g and, passing to a limit in the above, $Y(\omega) = g(X(\omega))$

Conversely, suppose $Y(\omega) = g(X(\omega))$. Why is Y $\sigma(X)$ measurable?

$$Y^{-1}(\text{open}) = X^{-1}(g^{-1}(\text{open})) = X^{-1}(\text{Borel}) \in \sigma(X) \blacksquare$$

26.5 Banach Space Valued Random Variables

Recall that for X a random variable, $\sigma(X)$ is the smallest σ algebra containing all the sets of the form $X^{-1}(F)$ where F is Borel. Since such sets, $X^{-1}(F)$ for F Borel form a σ algebra it follows $\sigma(X) = \{X^{-1}(F) : F \text{ is Borel}\}$.

Next consider the case where you have a set of σ algebras. The following lemma is helpful when you try to verify such a set of σ algebras is independent. It says you only need to check things on π systems contained in the σ algebras. This is really nice because it is much easier to consider the smaller π systems than the whole σ algebra.

Lemma 26.5.1 Suppose $\{\mathcal{F}_i\}_{i \in I}$ is a set of σ algebras contained in \mathcal{F} where \mathcal{F} is a σ algebra of sets of Ω . Suppose that $\mathcal{K}_i \subseteq \mathcal{F}_i$ is a π system and $\mathcal{F}_i = \sigma(\mathcal{K}_i)$. Suppose also that whenever J is a finite subset of I and $A_j \in \mathcal{K}_j$ for $j \in J$, it follows $P(\cap_{j \in J} A_j) = \prod_{j \in J} P(A_j)$. Then $\{\mathcal{F}_i\}_{i \in I}$ is independent.

Proof: I need to verify that under the given conditions, if $\{j_1, j_2, \dots, j_n\} \subseteq I$ and $A_{j_k} \subseteq \mathcal{F}_{j_k}$, then $P(\cap_{k=1}^n A_{j_k}) = \prod_{k=1}^n P(A_{j_k})$. By hypothesis, this is true if each $A_{j_k} \in \mathcal{K}_{j_k}$. Suppose it is true whenever there are at most $r-1 \geq 0$ of the A_{j_k} which are **not** in \mathcal{K}_{j_k} . Consider $\cap_{k=1}^n A_{j_k}$ where there are r sets which are **not** in the corresponding \mathcal{K}_{j_k} . Without loss of generality, say there are at most $r-1$ sets in the first $n-1$ which are not in the corresponding \mathcal{K}_{j_k} .

Pick $(A_{j_1}, \dots, A_{j_{n-1}})$ let

$$\mathcal{G}_{(A_{j_1}, \dots, A_{j_{n-1}})} \equiv \left\{ B \in \mathcal{F}_{j_n} : P(\cap_{k=1}^{n-1} A_{j_k} \cap B) = \prod_{k=1}^{n-1} P(A_{j_k}) P(B) \right\}$$

I am going to show $\mathcal{G}_{(A_{j_1}, \dots, A_{j_{n-1}})}$ is closed with respect to complements and countable disjoint unions and then apply the Lemma on π systems. By the induction hypothesis, $\mathcal{K}_{j_n} \subseteq \mathcal{G}_{(A_{j_1}, \dots, A_{j_{n-1}})}$. If $B \in \mathcal{G}_{(A_{j_1}, \dots, A_{j_{n-1}})}$,

$$\prod_{k=1}^{n-1} P(A_{j_k}) = P(\cap_{k=1}^{n-1} A_{j_k}) = P((\cap_{k=1}^{n-1} A_{j_k} \cap B^c) \cup (\cap_{k=1}^{n-1} A_{j_k} \cap B))$$

$$= P\left(\cap_{k=1}^{n-1} A_{j_k} \cap B^C\right) + P\left(\cap_{k=1}^{n-1} A_{j_k} \cap B\right) = P\left(\cap_{k=1}^{n-1} A_{j_k} \cap B^C\right) + \prod_{k=1}^{n-1} P(A_{j_k}) P(B)$$

and so

$$P\left(\cap_{k=1}^{n-1} A_{j_k} \cap B^C\right) = \prod_{k=1}^{n-1} P(A_{j_k}) (1 - P(B)) = \prod_{k=1}^{n-1} P(A_{j_k}) P(B^C)$$

showing if $B \in \mathcal{G}_{(A_{j_1}, \dots, A_{j_{n-1}})}$, then so is B^C . It is clear that $\mathcal{G}_{(A_{j_1}, \dots, A_{j_{n-1}})}$ is closed with respect to disjoint unions also. Here is why. If $\{B_j\}_{j=1}^\infty$ are disjoint sets in $\mathcal{G}_{(A_{j_1}, \dots, A_{j_{n-1}})}$,

$$\begin{aligned} P\left(\cup_{i=1}^\infty B_i \cap \cap_{k=1}^{n-1} A_{j_k}\right) &= \sum_{i=1}^\infty P\left(B_i \cap \cap_{k=1}^{n-1} A_{j_k}\right) = \sum_{i=1}^\infty P(B_i) \prod_{k=1}^{n-1} P(A_{j_k}) \\ &= \prod_{k=1}^{n-1} P(A_{j_k}) \sum_{i=1}^\infty P(B_i) = \prod_{k=1}^{n-1} P(A_{j_k}) P\left(\cup_{i=1}^\infty B_i\right) \end{aligned}$$

Therefore, by the π system lemma, Lemma 9.3.2 $\mathcal{G}_{(A_{j_1}, \dots, A_{j_{n-1}})} = \mathcal{F}_{j_n}$. This proves the induction step in going from $r-1$ to r . ■

What is a useful π system for $\mathcal{B}(E)$, the Borel sets of E where E is a Banach space?

Recall the fundamental lemma used to prove the Pettis theorem. It was proved on Page 649.

Lemma 26.5.2 *Let E be a separable real Banach space. Sets of the form*

$$\{x \in E : x_i^*(x) \leq \alpha_i, i = 1, 2, \dots, m\}$$

where $x_i^* \in D'$, a dense subspace of the unit ball of E' and $\alpha_i \in [-\infty, \infty)$ are a π system, and denoting this π system by \mathcal{K} , it follows $\sigma(\mathcal{K}) = \mathcal{B}(E)$. The sets of \mathcal{K} are examples of “cylindrical” sets. The D' is that set for the proof of the Pettis theorem.

Proof: The sets described are obviously a π system. I want to show $\sigma(\mathcal{K})$ contains the closed balls because then $\sigma(\mathcal{K})$ contains the open balls and hence the open sets and the result will follow. Let D' be described in Lemma 24.1.7. As pointed out earlier it can be any dense subset of B' . Then

$$\begin{aligned} \{x \in E : \|x - a\| \leq r\} &= \left\{x \in E : \sup_{f \in D'} |f(x - a)| \leq r\right\} \\ &= \left\{x \in E : \sup_{f \in D'} |f(x) - f(a)| \leq r\right\} \\ &= \cap_{f \in D'} \{x \in E : f(a) - r \leq f(x) \leq f(a) + r\} \\ &= \cap_{f \in D'} \{x \in E : f(x) \leq f(a) + r \text{ and } (-f)(x) \leq r - f(a)\} \end{aligned}$$

which equals a countable intersection of sets of the given π system. Therefore, every closed ball is contained in $\sigma(\mathcal{K})$. It follows easily that every open ball is also contained in $\sigma(\mathcal{K})$ because

$$B(a, r) = \cup_{n=1}^\infty \overline{B\left(a, r - \frac{1}{n}\right)}.$$

Since the Banach space is separable, it is completely separable and so every open set is the countable union of balls. This shows the open sets are in $\sigma(\mathcal{K})$ and so $\sigma(\mathcal{K}) \supseteq \mathcal{B}(E)$. However, all the sets in the π system are closed hence Borel because they are inverse images of closed sets. Therefore, $\sigma(\mathcal{K}) \subseteq \mathcal{B}(E)$ and so $\sigma(\mathcal{K}) = \mathcal{B}(E)$. ■

As mentioned above, we can replace D' in the above with M , any dense subset of E' .

Observation 26.5.3 Denote by $C_{\alpha,n}$ the set $\{\beta \in \mathbb{R}^n : \beta_i \leq \alpha_i\}$. Also denote by g_n an element of M^n where M is a dense subset of E' with the understanding that $g_n : E \rightarrow \mathbb{R}^n$ according to the rule

$$g_n(x) \equiv (g_1(x), \dots, g_n(x)).$$

Then the sets in the above lemma can be written as $g_n^{-1}(C_{\alpha,n})$. In other words, sets of the form $g_n^{-1}(C_{\alpha,n})$ form a π system for $\mathcal{B}(E)$.

Next suppose you have some random variables having values in a separable Banach space, E , $\{X_i\}_{i \in I}$. How can you tell if they are independent? To show they are independent, you need to verify that

$$P\left(\bigcap_{k=1}^n X_{i_k}^{-1}(F_{i_k})\right) = \prod_{k=1}^n P\left(X_{i_k}^{-1}(F_{i_k})\right)$$

whenever the F_{i_k} are Borel sets in E . It is desirable to find a way to do this easily.

Lemma 26.5.4 Let \mathcal{K} be a π system of sets of E , a separable real Banach space and let (Ω, \mathcal{F}, P) be a probability space and $X : \Omega \rightarrow E$ be a random variable. Then

$$X^{-1}(\sigma(\mathcal{K})) = \sigma(X^{-1}(\mathcal{K}))$$

Proof: First note that $X^{-1}(\sigma(\mathcal{K}))$ is a σ algebra which contains $X^{-1}(\mathcal{K})$ and so it contains $\sigma(X^{-1}(\mathcal{K}))$. Thus $X^{-1}(\sigma(\mathcal{K})) \supseteq \sigma(X^{-1}(\mathcal{K}))$. Now let

$$\mathcal{G} \equiv \{A \in \sigma(\mathcal{K}) : X^{-1}(A) \in \sigma(X^{-1}(\mathcal{K}))\}$$

Then $\mathcal{G} \supseteq \mathcal{K}$. If $A \in \mathcal{G}$, then $X^{-1}(A) \in \sigma(X^{-1}(\mathcal{K}))$ and so

$$X^{-1}(A)^C = X^{-1}(A^C) \in \sigma(X^{-1}(\mathcal{K}))$$

because $\sigma(X^{-1}(\mathcal{K}))$ is a σ algebra. Hence $A^C \in \mathcal{G}$. Finally suppose $\{A_i\}$ is a sequence of disjoint sets of \mathcal{G} . Then

$$X^{-1}(\cup_{i=1}^{\infty} A_i) = \cup_{i=1}^{\infty} X^{-1}(A_i) \in \sigma(X^{-1}(\mathcal{K}))$$

again because $\sigma(X^{-1}(\mathcal{K}))$ is a σ algebra. It follows from Lemma 9.3.2 on Page 243 that $\mathcal{G} \supseteq \sigma(\mathcal{K})$ and this shows that whenever

$$A \in \sigma(\mathcal{K}), X^{-1}(A) \in \sigma(X^{-1}(\mathcal{K})).$$

Thus $X^{-1}(\sigma(\mathcal{K})) \subseteq \sigma(X^{-1}(\mathcal{K}))$. ■

With this lemma, here is the desired result about independent random variables. Essentially, you can reduce to the case of random vectors having values in \mathbb{R}^n .

26.6 Reduction to Finite Dimensions

Let E be a Banach space and let $\mathbf{g} \in (E')^n$. Then for $x \in E$, $\mathbf{g} \circ x$ is the vector in \mathbb{R}^n which equals $(g_1(x), g_2(x), \dots, g_n(x))$.

Theorem 26.6.1 *Let X_i be a random variable having values in E a real separable Banach space. The random variables $\{X_i\}_{i \in I}$ are independent if and whenever*

$$\{i_1, \dots, i_n\} \subseteq I,$$

m_{i_1}, \dots, m_{i_n} are positive integers, and $\mathbf{g}_{m_{i_1}}, \dots, \mathbf{g}_{m_{i_n}}$ are respectively in

$$(M)^{m_{i_1}}, \dots, (M)^{m_{i_n}}$$

for M a dense subspace of E' , $\{\mathbf{g}_{m_{i_j}} \circ X_{i_j}\}_{j=1}^n$ are independent random vectors having values in $\mathbb{R}^{m_{i_1}}, \dots, \mathbb{R}^{m_{i_n}}$ respectively.

Proof: It is necessary to show that the events $X_{i_j}^{-1}(B_{i_j})$ are independent events whenever B_{i_j} are Borel sets. By Lemma 26.5.1 and the above Lemma 26.5.2, it suffices to verify that the events

$$X_{i_j}^{-1}(\mathbf{g}_{m_{i_j}}^{-1}(C_{\vec{\alpha}, m_{i_j}})) = (\mathbf{g}_{m_{i_j}} \circ X_{i_j})^{-1}(C_{\vec{\alpha}, m_{i_j}})$$

are independent where $C_{\vec{\alpha}, m_{i_j}}$ are the cones described in Lemma 26.5.2. Thus

$$\vec{\alpha} = (\alpha_{k_1}, \dots, \alpha_{k_m}), C_{\vec{\alpha}, m_{i_j}} = \prod_{i=1}^{m_{i_j}} (-\infty, \alpha_{k_i}]$$

But this condition is implied when the finite dimensional valued random vectors $\mathbf{g}_{m_{i_j}} \circ X_{i_j}$ are independent. ■

The above assertion also goes the other way as you may want to show.

26.7 0, 1 Laws

I am following [55] for the proof of many of the following theorems. Recall the set of ω which are in infinitely many of the sets $\{A_n\}$ is $\bigcap_{n=1}^{\infty} \bigcup_{m=n}^{\infty} A_m$. This is in $\bigcap_{n=1}^{\infty} \bigcup_{m=n}^{\infty} A_m$ if and only if for every n there exists $m \geq n$ such that it is in A_m .

Theorem 26.7.1 *Suppose $A_n \in \mathcal{F}_n$ where the σ algebras $\{\mathcal{F}_n\}_{n=1}^{\infty}$ are independent. Suppose also that $\sum_{k=1}^{\infty} P(A_k) = \infty$. Then $P(\bigcap_{n=1}^{\infty} \bigcup_{m=n}^{\infty} A_m) = 1$.*

Proof: It suffices to verify that $P(\bigcup_{n=1}^{\infty} \bigcap_{m=n}^{\infty} A_m^C) = 0$ which can be accomplished by showing that $P(\bigcap_{m=n}^{\infty} A_m^C) = 0$ for each n . The sets $\{A_k^C\}$ satisfy $A_k^C \in \mathcal{F}_k$. Therefore, noting that $e^{-x} \geq 1 - x$,

$$\begin{aligned} P(\bigcap_{m=n}^{\infty} A_m^C) &= \lim_{N \rightarrow \infty} P(\bigcap_{m=n}^N A_m^C) = \lim_{N \rightarrow \infty} \prod_{m=n}^N P(A_m^C) \\ &= \lim_{N \rightarrow \infty} \prod_{m=n}^N (1 - P(A_m)) \leq \lim_{N \rightarrow \infty} \prod_{m=n}^N e^{-P(A_m)} \\ &= \lim_{N \rightarrow \infty} \exp\left(-\sum_{m=n}^N P(A_m)\right) = 0. \quad \blacksquare \end{aligned}$$

The Kolmogorov zero one law follows next. It has to do with something called a tail event.

Definition 26.7.2 Let $\{\mathcal{F}_n\}$ be a sequence of σ algebras. Then $\mathcal{T}_n \equiv \sigma(\cup_{k=n}^{\infty} \mathcal{F}_k)$ where this means the smallest σ algebra which contains each \mathcal{F}_k for $k \geq n$. Then a tail event is a set which is in the σ algebra, $\mathcal{T} \equiv \cap_{n=1}^{\infty} \mathcal{T}_n$.

As usual, (Ω, \mathcal{F}, P) is the underlying probability space such that all σ algebras are contained in \mathcal{F} .

Lemma 26.7.3 Suppose $\{\mathcal{F}_n\}_{n=1}^{\infty}$ are independent σ algebras and suppose A is a tail event and $A_{k_i} \in \mathcal{F}_{k_i}$, $i = 1, \dots, m$ are given sets. Then

$$P(A_{k_1} \cap \dots \cap A_{k_m} \cap A) = P(A_{k_1} \cap \dots \cap A_{k_m}) P(A)$$

Proof: Let \mathcal{H} be the π system consisting of finite intersections of the form

$$B_{m_1} \cap B_{m_2} \cap \dots \cap B_{m_j}$$

where $B_{m_i} \in \mathcal{F}_{k_i}$ for $k_i > \max\{k_1, \dots, k_m\} \equiv N$. Thus $\sigma(\mathcal{H}) = \sigma(\cup_{i=N+1}^{\infty} \mathcal{F}_i)$. Now let

$$\mathcal{G} \equiv \{B \in \sigma(\mathcal{H}) : P(A_{k_1} \cap \dots \cap A_{k_m} \cap B) = P(A_{k_1} \cap \dots \cap A_{k_m}) P(B)\}$$

Then clearly $\mathcal{H} \subseteq \mathcal{G}$. It is also true that \mathcal{G} is closed with respect to complements and countable disjoint unions. By the lemma on π systems, $\mathcal{G} = \sigma(\mathcal{H}) = \sigma(\cup_{i=N+1}^{\infty} \mathcal{F}_i)$. Since A is in $\sigma(\cup_{i=N+1}^{\infty} \mathcal{F}_i)$ due to the assumption that it is a tail event, it follows that

$$P(A_{k_1} \cap \dots \cap A_{k_m} \cap A) = P(A_{k_1} \cap \dots \cap A_{k_m}) P(A) \blacksquare$$

Theorem 26.7.4 Suppose the σ algebras, $\{\mathcal{F}_n\}_{n=1}^{\infty}$ are independent and suppose A is a tail event. Then $P(A)$ either equals 0 or 1.

Proof: Let $A \in \mathcal{T} \equiv \cap_{n=1}^{\infty} \mathcal{T}_n \equiv \cap_{n=1}^{\infty} \sigma(\cup_{k=n}^{\infty} \mathcal{F}_k)$. I want to show that $P(A) = P(A)^2$. Since A is in \mathcal{T} , it is in each $\sigma(\cup_{k=n}^{\infty} \mathcal{F}_k)$. Let \mathcal{H} denote sets of the form $A_{k_1} \cap \dots \cap A_{k_m}$ for some m , $A_{k_j} \in \mathcal{F}_{k_j}$ where each $k_j > n$. Thus \mathcal{H} is a π system and

$$\sigma(\mathcal{H}) = \sigma(\cup_{k=n+1}^{\infty} \mathcal{F}_k) \equiv \mathcal{T}_{n+1}$$

Let

$$\mathcal{G} \equiv \{B \in \mathcal{T}_{n+1} \equiv \sigma(\cup_{k=n+1}^{\infty} \mathcal{F}_k) : P(A \cap B) = P(A) P(B)\}$$

Thus $\mathcal{H} \subseteq \mathcal{G}$ because

$$P(A_{k_1} \cap \dots \cap A_{k_m} \cap A) = P(A_{k_1} \cap \dots \cap A_{k_m}) P(A)$$

by Lemma 26.7.3. However, it is routine that \mathcal{G} is closed with respect to countable disjoint unions and complements. Therefore by the Lemma on π systems Lemma 9.3.2 on Page 243, it follows $\mathcal{G} = \sigma(\mathcal{H}) = \sigma(\cup_{k=n+1}^{\infty} \mathcal{F}_k)$.

Thus for any $B \in \sigma(\cup_{k=n+1}^{\infty} \mathcal{F}_k) = \mathcal{T}_{n+1}$, $P(A \cap B) = P(A) P(B)$. However, A is in all of these \mathcal{T}_{n+1} and so $P(A \cap A) = P(A) = P(A)^2$ so $P(A)$ equals either 0 or 1. \blacksquare

What sorts of things are tail events of independent σ algebras?

Theorem 26.7.5 *Let $\{X_k\}$ be a sequence of independent random variables having values in Z a Banach space. That is, the σ algebras $\sigma(X_k)$ are independent. Then*

$$A \equiv \{\omega : \{X_k(\omega)\} \text{ converges}\}$$

is a tail event. So is

$$B \equiv \left\{ \omega : \left\{ \sum_{k=1}^{\infty} X_k(\omega) \right\} \text{ converges} \right\}.$$

Proof: Since Z is complete, A is the same as the set where $\{X_k(\omega)\}$ is a Cauchy sequence. This set is

$$\bigcap_{n=1}^{\infty} \bigcap_{p=1}^{\infty} \bigcup_{m=p}^{\infty} \bigcap_{l,k \geq m} \{\omega : \|X_k(\omega) - X_l(\omega)\| < 1/n\}$$

Note that

$$\bigcup_{m=p}^{\infty} \bigcap_{l,k \geq m} \{\omega : \|X_k(\omega) - X_l(\omega)\| < 1/n\} \in \sigma(\bigcup_{j=p}^{\infty} \sigma(X_j))$$

for every p is the set where ultimately any pair of X_k, X_l are closer together than $1/n$,

$$\bigcap_{p=1}^{\infty} \bigcup_{m=p}^{\infty} \bigcap_{l,k \geq m} \{\omega : \|X_k(\omega) - X_l(\omega)\| < 1/n\}$$

is a tail event. The set where $\{X_k(\omega)\}$ is a Cauchy sequence is the intersection of all these and is therefore, also a tail event.

Now consider B . This set is the same as the set where the partial sums are Cauchy sequences. Let $S_n \equiv \sum_{k=1}^n X_k$. The set where the sum converges is then

$$\bigcap_{n=1}^{\infty} \bigcap_{p=2}^{\infty} \bigcup_{m=p}^{\infty} \bigcap_{l,k \geq m} \{\omega : \|S_k(\omega) - S_l(\omega)\| < 1/n\}$$

Say $k < l$ and consider for $m \geq p$

$$\{\omega : \|S_k(\omega) - S_l(\omega)\| < 1/n, k \geq m\}$$

This is the same as

$$\left\{ \omega : \left\| \sum_{j=k-1}^l X_j(\omega) \right\| < 1/n, k \geq m \right\} \in \sigma(\bigcup_{j=p-1}^{\infty} \sigma(X_j))$$

Thus

$$\bigcup_{m=p}^{\infty} \bigcap_{l,k \geq m} \{\omega : \|S_k(\omega) - S_l(\omega)\| < 1/n\} \in \sigma(\bigcup_{j=p-1}^{\infty} \sigma(X_j))$$

and so the intersection for all p of these is a tail event. Then the intersection over all n of these tail events is a tail event. ■

From this it can be concluded that if you have a sequence of independent random variables, $\{X_k\}$ the set where it converges is either of probability 1 or probability 0. A similar conclusion holds for the set where the infinite sum of these random variables converges. This is stated in the next corollary. This incredible assertion is the next corollary.

Corollary 26.7.6 *Let $\{X_k\}$ be a sequence of random variables having values in a Banach space. Then $\lim_{n \rightarrow \infty} X_n(\omega)$ either exists for a.e. ω or the convergence fails to take place for a.e. ω . Also if*

$$A \equiv \left\{ \omega : \sum_{k=1}^{\infty} X_k(\omega) \text{ converges} \right\},$$

then $P(A) = 0$ or 1 .

26.8 Strong Law of Large Numbers

Kolmogorov's inequality is a very interesting inequality which depends on independence of a set of random vectors. The random vectors have values in \mathbb{R}^n or more generally some real separable Hilbert space.

Lemma 26.8.1 *If Y, X are independent random variables having values in a real separable Hilbert space, H with $E(|X|^2), E(|Y|^2) < \infty$, then*

$$\int_{\Omega} (X, Y) dP = \left(\int_{\Omega} X dP, \int_{\Omega} Y dP \right).$$

Proof: Let $\{e_k\}$ be a complete orthonormal basis. Thus from Theorem 22.4.2,

$$\int_{\Omega} (X, Y) dP = \int_{\Omega} \sum_{k=1}^{\infty} (X, e_k) (Y, e_k) dP$$

Now

$$\begin{aligned} \int_{\Omega} \sum_{k=1}^{\infty} |(X, e_k) (Y, e_k)| dP &\leq \int_{\Omega} \left(\sum_k |(X, e_k)|^2 \right)^{1/2} \left(\sum_k |(Y, e_k)|^2 \right)^{1/2} dP \\ &= \int_{\Omega} |X| |Y| dP \leq \left(\int_{\Omega} |X|^2 dP \right)^{1/2} \left(\int_{\Omega} |Y|^2 dP \right)^{1/2} < \infty \end{aligned}$$

and so by Fubini's theorem and independence of X, Y ,

$$\begin{aligned} \int_{\Omega} (X, Y) dP &= \int_{\Omega} \sum_{k=1}^{\infty} (X, e_k) (Y, e_k) dP = \sum_{k=1}^{\infty} \int_{\Omega} (X, e_k) (Y, e_k) dP \\ &= \sum_{k=1}^{\infty} \int_{\Omega} (X, e_k) dP \int_{\Omega} (Y, e_k) dP = \sum_{k=1}^{\infty} \left(\int_{\Omega} X dP, e_k \right) \left(\int_{\Omega} Y dP, e_k \right) dP \\ &= \left(\int_{\Omega} X dP, \int_{\Omega} Y dP \right) \blacksquare \end{aligned}$$

Now here is Kolmogorov's inequality.

Theorem 26.8.2 *Suppose $\{X_k\}_{k=1}^n$ are independent with $E(|X_k|) < \infty$, $E(X_k) = 0$. Then for any $\varepsilon > 0$,*

$$P \left(\left[\max_{1 \leq k \leq n} \left| \sum_{j=1}^k X_j \right| \geq \varepsilon \right] \right) \leq \frac{1}{\varepsilon^2} \sum_{j=1}^n E(|X_j|^2).$$

Proof: Let $A = \left[\max_{1 \leq k \leq n} \left| \sum_{j=1}^k X_j \right| \geq \varepsilon \right]$. Now let $A_1 \equiv [|\sum_{j=1}^1 X_j| \geq \varepsilon]$ and if A_1, \dots, A_m have been chosen,

$$A_{m+1} \equiv \left[\left| \sum_{j=1}^{m+1} X_j \right| \geq \varepsilon \right] \cap \bigcap_{r=1}^m \left[\left| \sum_{j=1}^r X_j \right| < \varepsilon \right]$$

Thus the A_k partition A and $\omega \in A_k$ means $\left| \sum_{j=1}^k \mathbf{X}_j \right| \geq \varepsilon$ but this did not happen for $\left| \sum_{j=1}^r \mathbf{X}_j \right|$ for any $r < k$. Note also that $A_k \in \sigma(\mathbf{X}_1, \dots, \mathbf{X}_k)$. Then from algebra,

$$\begin{aligned} \left| \sum_{j=1}^n \mathbf{X}_j \right|^2 &= \left(\sum_{i=1}^k \mathbf{X}_i + \sum_{j=k+1}^n \mathbf{X}_j, \sum_{i=1}^k \mathbf{X}_i + \sum_{j=k+1}^n \mathbf{X}_j \right) \\ &= \left| \sum_{j=1}^k \mathbf{X}_j \right|^2 + \sum_{i \leq k, j > k} (\mathbf{X}_i, \mathbf{X}_j) + \sum_{i \leq k, j > k} (\mathbf{X}_j, \mathbf{X}_i) + \sum_{i > k, j > k} (\mathbf{X}_j, \mathbf{X}_i) \end{aligned}$$

Written more succinctly, $\left| \sum_{j=1}^n \mathbf{X}_j \right|^2 = \left| \sum_{j=1}^k \mathbf{X}_j \right|^2 + \sum_{j > k \text{ or } i > k} (\mathbf{X}_i, \mathbf{X}_j)$. Now multiply both sides by \mathcal{X}_{A_k} and integrate. Suppose $i \leq k$ for one of the terms in the second sum. Then by Lemma 26.4.4 and $A_k \in \sigma(\mathbf{X}_1, \dots, \mathbf{X}_k)$, the two random vectors $\mathcal{X}_{A_k} \mathbf{X}_i, \mathbf{X}_j$ are independent,

$$\int_{\Omega} \mathcal{X}_{A_k} (\mathbf{X}_i, \mathbf{X}_j) dP = \left(\int_{\Omega} \mathcal{X}_{A_k} \mathbf{X}_i dP, \int_{\Omega} \mathbf{X}_j dP \right) = 0$$

the last equality holding because by assumption $E(\mathbf{X}_j) = \mathbf{0}$. Therefore, it can be assumed both i, j are larger than k and

$$\int_{\Omega} \mathcal{X}_{A_k} \left| \sum_{j=1}^n \mathbf{X}_j \right|^2 dP = \int_{\Omega} \mathcal{X}_{A_k} \left| \sum_{j=1}^k \mathbf{X}_j \right|^2 dP + \sum_{j > k, i > k} \int_{\Omega} \mathcal{X}_{A_k} (\mathbf{X}_i, \mathbf{X}_j) dP \quad (26.8)$$

The last term on the right is interesting. Suppose $i > j$. The integral inside the sum is of the form $\int_{\Omega} (\mathbf{X}_i, \mathcal{X}_{A_k} \mathbf{X}_j) dP$. The second factor in the inner product is in

$$\sigma(\mathbf{X}_1, \dots, \mathbf{X}_k, \mathbf{X}_j)$$

and \mathbf{X}_i is not included in the list of random vectors. Thus by Lemma 26.4.4, the two random vectors $\mathbf{X}_i, \mathcal{X}_{A_k} \mathbf{X}_j$ are independent and so the last term in 26.8 reduces to

$$\left(\int_{\Omega} \mathbf{X}_i dP, \int_{\Omega} \mathcal{X}_{A_k} \mathbf{X}_j dP \right) = \left(\mathbf{0}, \int_{\Omega} \mathcal{X}_{A_k} \mathbf{X}_j dP \right) = 0.$$

A similar result holds if $j > i$. Thus the mixed terms in the last term of 26.8 are all equal to 0. Hence 26.8 reduces to

$$\int_{\Omega} \mathcal{X}_{A_k} \left| \sum_{j=1}^n \mathbf{X}_j \right|^2 dP = \int_{\Omega} \mathcal{X}_{A_k} \left| \sum_{j=1}^k \mathbf{X}_j \right|^2 dP + \sum_{i > k} \int_{\Omega} \mathcal{X}_{A_k} |\mathbf{X}_i|^2 dP$$

and so $\int_{\Omega} \mathcal{X}_{A_k} \left| \sum_{j=1}^n \mathbf{X}_j \right|^2 dP \geq \int_{\Omega} \mathcal{X}_{A_k} \left| \sum_{j=1}^k \mathbf{X}_j \right|^2 dP \geq \varepsilon^2 P(A_k)$. Now, summing these yields

$$\varepsilon^2 P(A) \leq \int_{\Omega} \mathcal{X}_A \left| \sum_{j=1}^n \mathbf{X}_j \right|^2 dP \leq \int_{\Omega} \left| \sum_{j=1}^n \mathbf{X}_j \right|^2 dP = \sum_{i,j} \int_{\Omega} (\mathbf{X}_i, \mathbf{X}_j) dP$$

By independence of the random vectors the mixed terms of the above sum equal zero and so it reduces to $\sum_{i=1}^n \int_{\Omega} |\mathbf{X}_i|^2 dP$ ■

This theorem implies the following amazing result.

Theorem 26.8.3 Let $\{X_k\}_{k=1}^\infty$ be independent random vectors having values in a separable real Hilbert space and suppose $E(|X_k|) < \infty$ for each k and $E(X_k) = \mathbf{0}$. Suppose also that $\sum_{j=1}^\infty E(|X_j|^2) < \infty$. Then $\sum_{j=1}^\infty X_j$ converges a.e.

Proof: Let $\varepsilon > 0$ be given. By Kolmogorov's inequality, Theorem 26.8.2, it follows that for $p \leq m < n$

$$P\left(\left[\max_{m \leq k \leq n} \left|\sum_{j=m}^k X_j\right| \geq \varepsilon\right]\right) \leq \frac{1}{\varepsilon^2} \sum_{j=p}^n E(|X_j|^2) \leq \frac{1}{\varepsilon^2} \sum_{j=p}^\infty E(|X_j|^2).$$

Therefore, letting $n \rightarrow \infty$ it follows that for all m, n such that $p \leq m \leq n$

$$P\left(\left[\max_{p \leq m \leq n} \left|\sum_{j=m}^n X_j\right| \geq \varepsilon\right]\right) \leq \frac{1}{\varepsilon^2} \sum_{j=p}^\infty E(|X_j|^2).$$

It follows from the assumption $\sum_{j=1}^\infty E(|X_j|^2) < \infty$ there exists a sequence, $\{p_n\}$ such that if $m \geq p_n$

$$P\left(\left[\max_{k \geq m \geq p_n} \left|\sum_{j=m}^k X_j\right| \geq 2^{-n}\right]\right) \leq 2^{-n}.$$

By the Borel Cantelli lemma, Lemma 26.1.2, there is a set of measure 0, N such that for $\omega \notin N$, ω is in only finitely many of the sets,

$$\left[\max_{k \geq m \geq p_n} \left|\sum_{j=m}^k X_j\right| \geq 2^{-n}\right]$$

and so for $\omega \notin N$, it follows that for large enough n ,

$$\left[\max_{k \geq m \geq p_n} \left|\sum_{j=m}^k X_j(\omega)\right| < 2^{-n}\right].$$

However, this says the partial sums $\left\{\sum_{j=1}^k X_j(\omega)\right\}_{k=1}^\infty$ are a Cauchy sequence. Therefore, they converge. ■

With this amazing result, there is a simple proof of the strong law of large numbers but first is an elementary lemma. In the following lemma, s_k and a_j could have values in any normed linear space.

Lemma 26.8.4 Suppose $s_k \rightarrow s$. Then $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n s_k = s$. Also if $\sum_{j=1}^\infty \frac{a_j}{j}$ converges, then $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n a_j = 0$.

Proof: Consider the first part. Since $s_k \rightarrow s$, it follows there is some constant, C such that $|s_k| < C$ for all k and $|s| < C$ also. Choose K so large that if $k \geq K$, then for $n > K$,

$$|s - s_k| < \varepsilon/2.$$

$$\left|s - \frac{1}{n} \sum_{k=1}^n s_k\right| \leq \frac{1}{n} \sum_{k=1}^n |s_k - s| = \frac{1}{n} \sum_{k=1}^K |s_k - s| + \frac{1}{n} \sum_{k=K}^n |s_k - s|$$

$$\leq \frac{2CK}{n} + \frac{\varepsilon}{2} \frac{n-K}{n} < \frac{2CK}{n} + \frac{\varepsilon}{2}$$

Therefore, whenever n is large enough, $|s - \frac{1}{n} \sum_{k=1}^n s_k| < \varepsilon$.

Now consider the second claim. Let $s_k = \sum_{j=1}^k \frac{a_j}{j}$ and $s = \lim_{k \rightarrow \infty} s_k$. Then by the first part,

$$\begin{aligned} s &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n s_k = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \sum_{j=1}^k \frac{a_j}{j} \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n \frac{a_j}{j} \sum_{k=j}^n 1 = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n \frac{a_j}{j} (n-j) \\ &= \lim_{n \rightarrow \infty} \left(\sum_{j=1}^n \frac{a_j}{j} - \frac{1}{n} \sum_{j=1}^n a_j \right) = s - \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n a_j \quad \blacksquare \end{aligned}$$

Now here is the strong law of large numbers.

Theorem 26.8.5 Suppose $\{X_k\}$ are independent random variables, and also suppose that $E(|X_k|) < \infty$ for each k and $E(X_k) = m_k$. Suppose also

$$\sum_{j=1}^{\infty} \frac{1}{j^2} E(|X_j - m_j|^2) < \infty. \quad (26.9)$$

Then $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n (X_j - m_j) = 0$.

Proof: Consider the sum $\sum_{j=1}^{\infty} \frac{X_j - m_j}{j}$. This sum converges *a.e.* because of 26.9 and Theorem 26.8.3 applied to the random vectors $\left\{ \frac{X_j - m_j}{j} \right\}$. Therefore, from Lemma 26.8.4 it follows that for a.e. ω , $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n (X_j(\omega) - m_j) = 0$ ■

The next corollary is often called the strong law of large numbers. It follows immediately from the above theorem.

Corollary 26.8.6 Suppose $\{X_j\}_{j=1}^{\infty}$ are independent having mean m and variance equal to $\sigma^2 \equiv \int_{\Omega} |X_j - m|^2 dP < \infty$. Then for a.e. $\omega \in \Omega$,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n X_j(\omega) = m$$

Chapter 27

Analytical Considerations

27.1 The Characteristic Function

One of the most important tools in probability is the characteristic function. To begin with, assume the random variables have values in \mathbb{R}^p .

Definition 27.1.1 Let X be a random variable as above. The characteristic function is

$$\phi_X(t) \equiv E(e^{it \cdot X}) \equiv \int_{\Omega} e^{it \cdot X(\omega)} dP = \int_{\mathbb{R}^p} e^{it \cdot x} d\lambda_X$$

the last equation holding by Proposition 26.1.12.

Recall the following fundamental lemma and definition, Lemma 13.2.4 on Page 379 where \mathcal{G} was a set of functions described there.

Definition 27.1.2 For $T \in \mathcal{G}^*$, define $FT, F^{-1}T \in \mathcal{G}^*$ by

$$FT(\phi) \equiv T(F\phi), F^{-1}T(\phi) \equiv T(F^{-1}\phi)$$

Lemma 27.1.3 F and F^{-1} are both one to one, onto, and are inverses of each other.

The main result on characteristic functions is the following.

Theorem 27.1.4 Let X and Y be random vectors with values in \mathbb{R}^p and suppose $E(e^{it \cdot X}) = E(e^{it \cdot Y})$ for all $t \in \mathbb{R}^p$. Then $\lambda_X = \lambda_Y$.

Proof: For $\psi \in \mathcal{G}$, let $\lambda_X(\psi) \equiv \int_{\mathbb{R}^p} \psi d\lambda_X$ and $\lambda_Y(\psi) \equiv \int_{\mathbb{R}^p} \psi d\lambda_Y$. Thus both λ_X and λ_Y are in \mathcal{G}^* . Then letting $\psi \in \mathcal{G}$ and using Fubini's theorem,

$$\begin{aligned} \int_{\mathbb{R}^p} \int_{\mathbb{R}^p} e^{it \cdot y} \psi(t) dt d\lambda_Y &= \int_{\mathbb{R}^p} \int_{\mathbb{R}^p} e^{it \cdot y} d\lambda_Y \psi(t) dt = \int_{\mathbb{R}^p} E(e^{it \cdot Y}) \psi(t) dt \\ &= \int_{\mathbb{R}^p} E(e^{it \cdot X}) \psi(t) dt = \int_{\mathbb{R}^p} \int_{\mathbb{R}^p} e^{it \cdot x} d\lambda_X \psi(t) dt \\ &= \int_{\mathbb{R}^p} \int_{\mathbb{R}^p} e^{it \cdot x} \psi(t) dt d\lambda_X. \end{aligned}$$

Thus $\lambda_Y(F^{-1}\psi) = \lambda_X(F^{-1}\psi)$. Since $\psi \in \mathcal{G}$ is arbitrary and F^{-1} is onto, this implies $\lambda_X = \lambda_Y$ in \mathcal{G}^* . But \mathcal{G} is dense in $C_0(\mathbb{R}^p)$ from the Stone Weierstrass theorem and so $\lambda_X = \lambda_Y$ as measures. Recall from real analysis the dual space of $C_0(\mathbb{R}^n)$ is the space of complex measures.

Alternatively, the above shows that since F^{-1} is onto, for all $\psi \in \mathcal{G}$,

$$\int_{\mathbb{R}^p} \psi d\lambda_Y = \int_{\mathbb{R}^p} \psi d\lambda_X$$

and then, by a use of the Stone Weierstrass theorem, the above will hold for all $\psi \in C_c(\mathbb{R}^n)$ and now, by the Riesz representation theorem for positive linear functionals, $\lambda_Y = \lambda_X$. ■

You can also give a version of this theorem in which reference is made only to the probability distribution measures.

Definition 27.1.5 For μ a probability measure on the Borel sets of \mathbb{R}^n ,

$$\phi_\mu(t) \equiv \int_{\mathbb{R}^n} e^{it \cdot x} d\mu.$$

Theorem 27.1.6 Let μ and ν be probability measures on the Borel sets of \mathbb{R}^p and suppose $\phi_\mu(t) = \phi_\nu(t)$. Then $\mu = \nu$.

Proof: The proof is identical to the above. Just replace λ_X with μ and λ_Y with ν . ■

27.2 Conditional Probability

Here I will consider the concept of conditional probability depending on the theory of differentiation of general Radon measures. This leads to a different way of thinking about independence.

If X, Y are random vectors defined on a probability space having values in \mathbb{R}^{p_1} and \mathbb{R}^{p_2} respectively, and if E is a Borel set in the appropriate space, then (X, Y) is a random vector with values in $\mathbb{R}^{p_1} \times \mathbb{R}^{p_2}$ and $\lambda_{(X,Y)}(E \times \mathbb{R}^{p_2}) = \lambda_X(E)$, $\lambda_{(X,Y)}(\mathbb{R}^{p_1} \times E) = \lambda_Y(E)$. Thus, by Theorem 19.8.1 on Page 520, there exist probability measures, denoted here by $\lambda_{X|Y}$ and $\lambda_{Y|X}$, such that whenever E is a Borel set in $\mathbb{R}^{p_1} \times \mathbb{R}^{p_2}$,

$$\int_{\mathbb{R}^{p_1} \times \mathbb{R}^{p_2}} \mathcal{X}_E d\lambda_{(X,Y)} = \int_{\mathbb{R}^{p_1}} \int_{\mathbb{R}^{p_2}} \mathcal{X}_E d\lambda_{Y|X} d\lambda_X,$$

and

$$\int_{\mathbb{R}^{p_1} \times \mathbb{R}^{p_2}} \mathcal{X}_E d\lambda_{(X,Y)} = \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_1}} \mathcal{X}_E d\lambda_{X|Y} d\lambda_Y.$$

Definition 27.2.1 Let X and Y be two random vectors defined on a probability space. The conditional probability measure of Y given X is the measure $\lambda_{Y|X}$ in the above. Similarly the conditional probability measure of X given Y is the measure $\lambda_{X|Y}$.

More generally, one can use the theory of slicing measures to consider any finite list of random vectors, $\{X_i\}$, defined on a probability space with $X_i \in \mathbb{R}^{p_i}$, and write the following for E a Borel set in $\prod_{i=1}^n \mathbb{R}^{p_i}$.

$$\begin{aligned} \int_{\mathbb{R}^{p_1} \times \dots \times \mathbb{R}^{p_n}} \mathcal{X}_E d\lambda_{(X_1, \dots, X_n)} &= \int_{\mathbb{R}^{p_1} \times \dots \times \mathbb{R}^{p_{n-1}}} \int_{\mathbb{R}^{p_n}} \mathcal{X}_E d\lambda_{X_n|(x_1, \dots, x_{n-1})} d\lambda_{(X_1, \dots, X_{n-1})} \\ &= \int_{\mathbb{R}^{p_1} \times \dots \times \mathbb{R}^{p_{n-2}}} \int_{\mathbb{R}^{p_{n-1}}} \int_{\mathbb{R}^{p_n}} \mathcal{X}_E d\lambda_{X_n|(x_1, \dots, x_{n-1})} d\lambda_{X_{n-1}|(x_1, \dots, x_{n-2})} d\lambda_{(X_1, \dots, X_{n-2})} \\ &\quad \vdots \\ &= \int_{\mathbb{R}^{p_1}} \dots \int_{\mathbb{R}^{p_n}} \mathcal{X}_E d\lambda_{X_n|(x_1, \dots, x_{n-1})} d\lambda_{X_{n-1}|(x_1, \dots, x_{n-2})} \dots d\lambda_{X_2|x_1} d\lambda_{X_1}. \end{aligned} \quad (27.1)$$

Obviously, this could have been done in any order in the iterated integrals by simply modifying the “given” variables, those occurring after the symbol $|$, to be those which have been integrated in an outer level of the iterated integral. For simplicity, write

$$\lambda_{X_n|(x_1, \dots, x_{n-1})} = \lambda_{X_n|x_1, \dots, x_{n-1}}$$

Definition 27.2.2 Let $\{\mathbf{X}_1, \dots, \mathbf{X}_n\}$ be random vectors defined on a probability space having values in $\mathbb{R}^{p_1}, \dots, \mathbb{R}^{p_n}$ respectively. The random vectors are independent if for every E a Borel set in $\mathbb{R}^{p_1} \times \dots \times \mathbb{R}^{p_n}$,

$$\begin{aligned} & \int_{\mathbb{R}^{p_1} \times \dots \times \mathbb{R}^{p_n}} \mathcal{X}_E d\lambda_{(\mathbf{X}_1, \dots, \mathbf{X}_n)} \\ &= \int_{\mathbb{R}^{p_1}} \dots \int_{\mathbb{R}^{p_n}} \mathcal{X}_E d\lambda_{\mathbf{X}_n} d\lambda_{\mathbf{X}_{n-1}} \dots d\lambda_{\mathbf{X}_2} d\lambda_{\mathbf{X}_1} \end{aligned} \quad (27.2)$$

and the iterated integration may be taken in any order. If \mathcal{A} is any set of random vectors defined on a probability space, \mathcal{A} is independent if any finite set of random vectors from \mathcal{A} is independent.

Thus, the random vectors are independent exactly when the dependence on the givens in 27.1 can be dropped.

Does this amount to the same thing as discussed earlier? Suppose you have three random variables $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$. Let $A = \mathbf{X}^{-1}(E)$, $B = \mathbf{Y}^{-1}(F)$, $C = \mathbf{Z}^{-1}(G)$ where E, F, G are Borel sets. Thus these inverse images are typical sets in

$$\sigma(\mathbf{X}), \sigma(\mathbf{Y}), \sigma(\mathbf{Z})$$

respectively. First suppose that the random variables are independent in the earlier sense. Then

$$\begin{aligned} P(A \cap B \cap C) &= P(A)P(B)P(C) \\ &= \int_{\mathbb{R}^{p_1}} \mathcal{X}_E(\mathbf{x}) d\lambda_{\mathbf{X}} \int_{\mathbb{R}^{p_2}} \mathcal{X}_F(\mathbf{y}) d\lambda_{\mathbf{Y}} \int_{\mathbb{R}^{p_3}} \mathcal{X}_G(\mathbf{z}) d\lambda_{\mathbf{Z}} \\ &= \int_{\mathbb{R}^{p_1}} \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_3}} \mathcal{X}_E(\mathbf{x}) \mathcal{X}_F(\mathbf{y}) \mathcal{X}_G(\mathbf{z}) d\lambda_{\mathbf{Z}} d\lambda_{\mathbf{Y}} d\lambda_{\mathbf{X}} \end{aligned}$$

Also

$$\begin{aligned} P(A \cap B \cap C) &= \int_{\mathbb{R}^{p_1} \times \mathbb{R}^{p_2} \times \mathbb{R}^{p_3}} \mathcal{X}_E(\mathbf{x}) \mathcal{X}_F(\mathbf{y}) \mathcal{X}_G(\mathbf{z}) d\lambda_{(\mathbf{X}, \mathbf{Y}, \mathbf{Z})} \\ &= \int_{\mathbb{R}^{p_1}} \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_3}} \mathcal{X}_E(\mathbf{x}) \mathcal{X}_F(\mathbf{y}) \mathcal{X}_G(\mathbf{z}) d\lambda_{\mathbf{Z}|\mathbf{x}\mathbf{y}} d\lambda_{\mathbf{Y}|\mathbf{x}} d\lambda_{\mathbf{X}} \end{aligned}$$

Thus

$$\begin{aligned} & \int_{\mathbb{R}^{p_1}} \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_3}} \mathcal{X}_E(\mathbf{x}) \mathcal{X}_F(\mathbf{y}) \mathcal{X}_G(\mathbf{z}) d\lambda_{\mathbf{Z}} d\lambda_{\mathbf{Y}} d\lambda_{\mathbf{X}} \\ &= \int_{\mathbb{R}^{p_1}} \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_3}} \mathcal{X}_E(\mathbf{x}) \mathcal{X}_F(\mathbf{y}) \mathcal{X}_G(\mathbf{z}) d\lambda_{\mathbf{Z}|\mathbf{x}\mathbf{y}} d\lambda_{\mathbf{Y}|\mathbf{x}} d\lambda_{\mathbf{X}} \end{aligned}$$

Now letting $G = \mathbb{R}^{p_3}$, it follows that

$$\int_{\mathbb{R}^{p_1}} \int_{\mathbb{R}^{p_2}} \mathcal{X}_E(\mathbf{x}) \mathcal{X}_F(\mathbf{y}) d\lambda_{\mathbf{Y}} d\lambda_{\mathbf{X}} = \int_{\mathbb{R}^{p_1}} \int_{\mathbb{R}^{p_2}} \mathcal{X}_E(\mathbf{x}) \mathcal{X}_F(\mathbf{y}) d\lambda_{\mathbf{Y}|\mathbf{x}} d\lambda_{\mathbf{X}}$$

By uniqueness of the slicing measures or an application of the Besikovitch differentiation theorem, it follows that for $\lambda_{\mathbf{X}}$ a.e. \mathbf{x} ,

$$\lambda_{\mathbf{Y}} = \lambda_{\mathbf{Y}|\mathbf{x}}$$

Thus, using this in the above,

$$\begin{aligned} & \int_{\mathbb{R}^{p_1}} \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_3}} \mathcal{X}_E(\mathbf{x}) \mathcal{X}_F(\mathbf{y}) \mathcal{X}_G(\mathbf{z}) d\lambda_{\mathbf{Z}} d\lambda_{\mathbf{Y}} d\lambda_{\mathbf{X}} \\ &= \int_{\mathbb{R}^{p_1}} \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_3}} \mathcal{X}_E(\mathbf{x}) \mathcal{X}_F(\mathbf{y}) \mathcal{X}_G(\mathbf{z}) d\lambda_{\mathbf{Z}|\mathbf{x}\mathbf{y}} d\lambda_{\mathbf{Y}} d\lambda_{\mathbf{X}} \end{aligned}$$

and also it reduces to

$$\begin{aligned} & \int_{\mathbb{R}^{p_1} \times \mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_3}} \mathcal{X}_E(\mathbf{x}) \mathcal{X}_F(\mathbf{y}) \mathcal{X}_G(\mathbf{z}) d\lambda_{\mathbf{Z}} d\lambda_{(\mathbf{X}, \mathbf{Y})} \\ &= \int_{\mathbb{R}^{p_1} \times \mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_3}} \mathcal{X}_E(\mathbf{x}) \mathcal{X}_F(\mathbf{y}) \mathcal{X}_G(\mathbf{z}) d\lambda_{\mathbf{Z}|\mathbf{x}\mathbf{y}} d\lambda_{(\mathbf{X}, \mathbf{Y})} \end{aligned}$$

Now by uniqueness of the slicing measures again, for $\lambda_{(\mathbf{X}, \mathbf{Y})}$ a.e. (\mathbf{x}, \mathbf{y}) , it follows that

$$\lambda_{\mathbf{Z}} = \lambda_{\mathbf{Z}|\mathbf{x}\mathbf{y}}$$

Similar conclusions hold for $\lambda_{\mathbf{X}}, \lambda_{\mathbf{Y}}$. In each case, off a set of measure zero the distribution measures equal the slicing measures.

Conversely, if the distribution measures equal the slicing measures off sets of measure zero as described above, then it is obvious that the random variables are independent. The same reasoning applies for any number of random variables.

Thus this gives a different and more analytical way to think of independence of finitely many random variables. Clearly, the argument given above will apply to any finite set of random variables.

Proposition 27.2.3 *Equations 27.2 and 27.1 hold with \mathcal{X}_E replaced by any nonnegative Borel measurable function and for any bounded continuous function or for any function in L^1 .*

Proof: The two equations hold for simple functions in place of \mathcal{X}_E and so an application of the monotone convergence theorem applied to an increasing sequence of simple functions converging pointwise to a given nonnegative Borel measurable function yields the conclusion of the proposition in the case of the nonnegative Borel function. For a bounded continuous function or one in L^1 , one can apply the result just established to the positive and negative parts of the real and imaginary parts of the function.

Lemma 27.2.4 *Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be random vectors with values in $\mathbb{R}^{p_1}, \dots, \mathbb{R}^{p_n}$ respectively and let*

$$\mathbf{g} : \mathbb{R}^{p_1} \times \dots \times \mathbb{R}^{p_n} \rightarrow \mathbb{R}^k$$

be Borel measurable. Then $\mathbf{g}(\mathbf{X}_1, \dots, \mathbf{X}_n)$ is a random vector with values in \mathbb{R}^k and if $h : \mathbb{R}^k \rightarrow [0, \infty)$, then

$$\begin{aligned} & \int_{\mathbb{R}^k} h(\mathbf{y}) d\lambda_{\mathbf{g}(\mathbf{X}_1, \dots, \mathbf{X}_n)}(\mathbf{y}) = \\ & \int_{\mathbb{R}^{p_1} \times \dots \times \mathbb{R}^{p_n}} h(\mathbf{g}(\mathbf{x}_1, \dots, \mathbf{x}_n)) d\lambda_{(\mathbf{X}_1, \dots, \mathbf{X}_n)}. \end{aligned} \quad (27.3)$$

If \mathbf{X}_i is a random vector with values in $\mathbb{R}^{p_i}, i = 1, 2, \dots$ and if $\mathbf{g}_i : \mathbb{R}^{p_i} \rightarrow \mathbb{R}^{k_i}$, where \mathbf{g}_i is Borel measurable, then the random vectors $\mathbf{g}_i(\mathbf{X}_i)$ are also independent whenever the \mathbf{X}_i are independent.

Proof: First let E be a Borel set in \mathbb{R}^k . From the definition,

$$\begin{aligned}\lambda_{g(\mathbf{X}_1, \dots, \mathbf{X}_n)}(E) &= P(g(\mathbf{X}_1, \dots, \mathbf{X}_n) \in E) \\ &= P((\mathbf{X}_1, \dots, \mathbf{X}_n) \in g^{-1}(E)) = \lambda_{(\mathbf{X}_1, \dots, \mathbf{X}_n)}(g^{-1}(E)) \\ \int_{\mathbb{R}^k} \mathcal{X}_E d\lambda_{g(\mathbf{X}_1, \dots, \mathbf{X}_n)} &= \int_{\mathbb{R}^{p_1} \times \dots \times \mathbb{R}^{p_n}} \mathcal{X}_{g^{-1}(E)} d\lambda_{(\mathbf{X}_1, \dots, \mathbf{X}_n)} \\ &= \int_{\mathbb{R}^{p_1} \times \dots \times \mathbb{R}^{p_n}} \mathcal{X}_E(g(\mathbf{x}_1, \dots, \mathbf{x}_n)) d\lambda_{(\mathbf{X}_1, \dots, \mathbf{X}_n)}.\end{aligned}$$

This proves 27.3 in the case when h is \mathcal{X}_E . To prove it in the general case, approximate the nonnegative Borel measurable function with simple functions for which the formula is true, and use the monotone convergence theorem.

It remains to prove the last assertion that functions of independent random vectors are also independent random vectors. Let E be a Borel set in $\mathbb{R}^{k_1} \times \dots \times \mathbb{R}^{k_n}$. Then for

$$\begin{aligned}\pi_i(\mathbf{x}_1, \dots, \mathbf{x}_n) &\equiv \mathbf{x}_i, \\ \int_{\mathbb{R}^{k_1} \times \dots \times \mathbb{R}^{k_n}} \mathcal{X}_E d\lambda_{(g_1(\mathbf{X}_1), \dots, g_n(\mathbf{X}_n))} \\ &\equiv \int_{\mathbb{R}^{p_1} \times \dots \times \mathbb{R}^{p_n}} \mathcal{X}_E \circ (g_1 \circ \pi_1, \dots, g_n \circ \pi_n) d\lambda_{(\mathbf{X}_1, \dots, \mathbf{X}_n)} \\ &= \int_{\mathbb{R}^{p_1}} \dots \int_{\mathbb{R}^{p_n}} \mathcal{X}_E \circ (g_1 \circ \pi_1, \dots, g_n \circ \pi_n) d\lambda_{\mathbf{X}_n} \dots d\lambda_{\mathbf{X}_1} \\ &= \int_{\mathbb{R}^{k_1}} \dots \int_{\mathbb{R}^{k_n}} \mathcal{X}_E d\lambda_{g_n(\mathbf{X}_n)} \dots d\lambda_{g_1(\mathbf{X}_1)} \blacksquare\end{aligned}$$

Of course if $X_i, i = 1, 2, \dots, n$ are independent, this means the σ algebras $\sigma(X_i)$ are independent. Now $\sigma(g_i \circ X_i) \subseteq \sigma(X_i)$ because

$$(g_i \circ X_i)^{-1}(\text{Borel set}) = X_i^{-1}(g_i^{-1}(\text{Borel set})) = X_i^{-1}(\text{Borel set}) \in \sigma(X_i)$$

and so the variables $g_i \circ X_i, i = 1, 2, \dots, n$ are independent. I think this is a more direct way of seeing this second claim.

Proposition 27.2.5 *Let ν_1, \dots, ν_n be Radon probability measures defined on \mathbb{R}^p . Then there exists a probability space and independent random vectors*

$$\{\mathbf{X}_1, \dots, \mathbf{X}_n\}$$

defined on this probability space such that $\lambda_{\mathbf{X}_i} = \nu_i$.

Proof: Let $(\Omega, \mathcal{S}, P) \equiv ((\mathbb{R}^p)^n, \mathcal{S}_1 \times \dots \times \mathcal{S}_n, \nu_1 \times \dots \times \nu_n)$ where this is just the product σ algebra and product measure which satisfies the following for measurable rectangles.

$$(\nu_1 \times \dots \times \nu_n) \left(\prod_{i=1}^n E_i \right) = \prod_{i=1}^n \nu_i(E_i).$$

Now let $\mathbf{X}_i(\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_n) = \mathbf{x}_i$.

Then from the definition, if E is a Borel set in \mathbb{R}^p ,

$$\begin{aligned}\lambda_{\mathbf{X}_i}(E) &\equiv P\{\mathbf{X}_i \in E\} \\ &= (\mathbf{v}_1 \times \cdots \times \mathbf{v}_n)(\mathbb{R}^p \times \cdots \times E \times \cdots \times \mathbb{R}^p) = \mathbf{v}_i(E).\end{aligned}$$

This defines the random vectors $\{\mathbf{X}_1, \dots, \mathbf{X}_n\}$ such that $\lambda_{\mathbf{X}_i} = \mathbf{v}_i$ on all Borel sets. Are these random vectors independent? Let \mathcal{M} consist of all Borel sets of $(\mathbb{R}^p)^n$ such that

$$\int_{\mathbb{R}^p} \cdots \int_{\mathbb{R}^p} \mathcal{X}_E(\mathbf{x}_1, \dots, \mathbf{x}_n) d\lambda_{\mathbf{X}_1} \cdots d\lambda_{\mathbf{X}_n} = \int_{(\mathbb{R}^p)^n} \mathcal{X}_E d\lambda_{(\mathbf{X}_1, \dots, \mathbf{X}_n)}.$$

It is clear that \mathcal{M} contains all products of Borel sets $\prod_{i=1}^n E_i$ which is a π system called \mathcal{K} . It is also clearly closed with respect to countable disjoint unions and complements. Thus, by the lemma on π systems, Lemma 9.3.2, it contains $\sigma(\mathcal{K})$ which is the Borel sets because it contains all open sets. Therefore, the given random vectors are independent because you can dispense with the givens. ■

The following Lemma was proved earlier in a different way.

Lemma 27.2.6 *If $\{X_i\}_{i=1}^n$ are independent random variables having values in \mathbb{R} ,*

$$E\left(\prod_{i=1}^n X_i\right) = \prod_{i=1}^n E(X_i).$$

Proof: By Lemma 27.2.4 and denoting by P the product, $\prod_{i=1}^n X_i$,

$$\begin{aligned}E\left(\prod_{i=1}^n X_i\right) &= \int_{\mathbb{R}} z d\lambda_P(z) = \int_{\mathbb{R}^n} \prod_{i=1}^n x_i d\lambda_{(X_1, \dots, X_n)} \\ &= \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} \prod_{i=1}^n x_i d\lambda_{X_1} \cdots d\lambda_{X_n} = \prod_{i=1}^n E(X_i). \quad \blacksquare\end{aligned}$$

27.3 Conditional Expectation, Sub-martingales

This concept is developed more later when conditional expectation relative to a σ algebra is discussed. However, I think there tends to be a disconnect between that more abstract idea and what we usually think of where conditional expectation involves expectation in which conditional probability is used, where a random variable is “given” to have a particular value, as described above. Certainly this is the case in beginning treatments of probability as an application of combinatorics or in beginning statistics classes. Therefore, to tie this in to this more elementary way of thinking, I will present conditional expectation in terms of conditional probability as defined above, where values of random variables are given. This leads to a rudimentary treatment of the sub-martingale convergence theorem presented here. However, the more abstract version presented later is much less difficult I think, because it is dependent on the Radon Nikodym theorem rather than the very difficult Besicovitch differentiation theory of Radon measures.

Definition 27.3.1 *Let X and Y be random vectors having values in \mathbb{F}^{p_1} and \mathbb{F}^{p_2} respectively. Then if $\int |x| d\lambda_{\mathbf{X}|\mathbf{Y}}(x) < \infty$, define $E(X|\mathbf{Y}) \equiv \int x d\lambda_{\mathbf{X}|\mathbf{Y}}(x)$.*

Proposition 27.3.2 Suppose $\int_{\mathbb{R}^{p_1} \times \mathbb{R}^{p_2}} |x| d\lambda_{(X,Y)}(x) < \infty$. Then $E(X|y)$ exists for λ_Y a.e. y and

$$\int_{\mathbb{R}^{p_2}} E(X|y) d\lambda_Y = \int_{\mathbb{R}^{p_1}} x d\lambda_X(x) = E(X).$$

Proof: $\infty > \int_{\mathbb{R}^{p_1} \times \mathbb{R}^{p_2}} |x| d\lambda_{(X,Y)} = \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_1}} |x| d\lambda_{X|y}(x) d\lambda_Y(y)$ and so

$$\int_{\mathbb{R}^{p_1}} |x| d\lambda_{X|y}(x) < \infty,$$

λ_Y a.e. Now $\int_{\mathbb{R}^{p_2}} E(X|y) d\lambda_Y =$

$$\begin{aligned} &= \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_1}} x d\lambda_{X|y}(x) d\lambda_Y(y) = \int_{\mathbb{R}^{p_1} \times \mathbb{R}^{p_2}} x d\lambda_{(X,Y)} \\ &= \int_{\mathbb{R}^{p_1}} \int_{\mathbb{R}^{p_2}} x d\lambda_{Y|x}(y) d\lambda_X(x) = \int_{\mathbb{R}^{p_1}} x d\lambda_X(x) = E(X) \blacksquare \end{aligned}$$

Definition 27.3.3 Let $\{X_n\}$ be any sequence, finite or infinite, of random variables with values in \mathbb{R} which are defined on some probability space, (Ω, \mathcal{S}, P) . We say $\{X_n\}$ is a martingale if

$$E(X_n | x_{n-1}, \dots, x_1) = x_{n-1}$$

and we say $\{X_n\}$ is a sub-martingale if

$$E(X_n | x_{n-1}, \dots, x_1) \geq x_{n-1}.$$

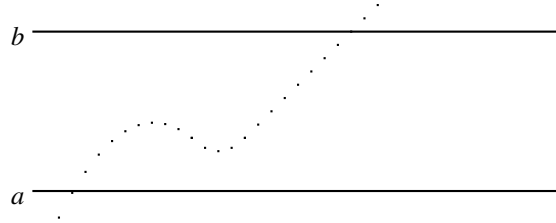
Recall Lemma 10.15.1, Jensen's inequality. It is stated next for convenience.

Lemma 27.3.4 If $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is convex, then ϕ is continuous. Also, if ϕ is convex, $\mu(\Omega) = 1$, and $f, \phi(f) : \Omega \rightarrow \mathbb{R}$ are in $L^1(\Omega)$, then $\phi(\int_{\Omega} f d\mu) \leq \int_{\Omega} \phi(f) d\mu$.

Next is the notion of an upcrossing.

Definition 27.3.5 Let $\{x_i\}_{i=1}^I$ be any sequence of real numbers, $I \leq \infty$. Define an increasing sequence of integers $\{m_k\}$ as follows. m_1 is the first integer ≥ 1 such that $x_{m_1} \leq a$, m_2 is the first integer larger than m_1 such that $x_{m_2} \geq b$, m_3 is the first integer larger than m_2 such that $x_{m_3} \leq a$, etc. Then each sequence, $\{x_{m_{2k-1}}, \dots, x_{m_{2k}}\}$, is called an upcrossing of $[a, b]$.

Here is a picture of an upcrossing.



Proposition 27.3.6 Let $\{X_i\}_{i=1}^n$ be a finite sequence of real random variables defined on Ω where (Ω, \mathcal{S}, P) is a probability space. Let $U_{[a,b]}(\omega)$ denote the number of upcrossings of $X_i(\omega)$ of the interval $[a, b]$. Then $U_{[a,b]}$ is a random variable, in other words, a nonnegative measurable function.

Proof: Let $X_0(\omega) \equiv a + 1$, let $Y_0(\omega) \equiv 0$, and let $Y_k(\omega)$ remain 0 for $k = 0, \dots, l$ until $X_l(\omega) \leq a$. When this happens (if ever), $Y_{l+1}(\omega) \equiv 1$. Then let $Y_i(\omega)$ remain 1 for $i = l + 1, \dots, r$ until $X_r(\omega) \geq b$ when $Y_{r+1}(\omega) \equiv 0$. Let $Y_k(\omega)$ remain 0 for $k \geq r + 1$ until $X_k(\omega) \leq a$ when $Y_k(\omega) \equiv 1$ and continue in this way. Thus the upcrossings of $X_i(\omega)$ are identified as unbroken strings of ones with a zero at each end, with the possible exception of the last string of ones which may be missing the zero at the upper end and may or may not be an upcrossing.

Note also that Y_0 is measurable because it is identically equal to 0 and that if Y_k is measurable, then Y_{k+1} is measurable because the only change in going from k to $k + 1$ is a change from 0 to 1 or from 1 to 0 on a measurable set determined by X_k . Now let

$$Z_k(\omega) = \begin{cases} 1 & \text{if } Y_k(\omega) = 1 \text{ and } Y_{k+1}(\omega) = 0, \\ 0 & \text{otherwise,} \end{cases}$$

if $k < n$ and

$$Z_n(\omega) = \begin{cases} 1 & \text{if } Y_n(\omega) = 1 \text{ and } X_n(\omega) \geq b, \\ 0 & \text{otherwise.} \end{cases}$$

Thus $Z_k(\omega) = 1$ exactly when an upcrossing has been completed and each Z_i is a random variable.

$$U_{[a,b]}(\omega) = \sum_{k=1}^n Z_k(\omega)$$

so $U_{[a,b]}$ is a random variable as claimed. ■

The following corollary collects some key observations found in the above construction.

Corollary 27.3.7 $U_{[a,b]}(\omega) \leq$ the number of unbroken strings of ones in the sequence $\{Y_k(\omega)\}$ there being at most one unbroken string of ones which produces no upcrossing. Also

$$Y_i(\omega) = \psi_i\left(\{X_j(\omega)\}_{j=1}^{i-1}\right), \quad (27.4)$$

where ψ_i is some function of the past values of $X_j(\omega)$.

Lemma 27.3.8 (upcrossing lemma) Let $\{X_i\}_{i=1}^n$ be a sub-martingale and suppose

$$E(|X_n|) < \infty.$$

Then

$$E(U_{[a,b]}) \leq \frac{E(|X_n|) + |a|}{b - a}.$$

Proof: Let $\phi(x) \equiv a + (x - a)^+$. Thus ϕ is a convex and increasing function.

$$\begin{aligned} \phi(X_{k+r}) - \phi(X_k) &= \sum_{i=k+1}^{k+r} \phi(X_i) - \phi(X_{i-1}) \\ &= \sum_{i=k+1}^{k+r} (\phi(X_i) - \phi(X_{i-1})) Y_i + \sum_{i=k+1}^{k+r} (\phi(X_i) - \phi(X_{i-1})) (1 - Y_i). \end{aligned}$$

The upcrossings of $\phi(X_i)$ are exactly the same as the upcrossings of X_i and from 27.4,

$$E\left(\sum_{i=k+1}^{k+r} (\phi(X_i) - \phi(X_{i-1})) (1 - Y_i)\right)$$

$$\begin{aligned}
&= \sum_{i=k+1}^{k+r} \int_{\mathbb{R}^i} (\phi(x_i) - \phi(x_{i-1})) \left(1 - \psi_i\left(\{x_j\}_{j=1}^{i-1}\right)\right) d\lambda_{(X_1, \dots, X_i)} \\
&= \sum_{i=k+1}^{k+r} \int_{\mathbb{R}^{i-1}} \int_{\mathbb{R}} (\phi(x_i) - \phi(x_{i-1})) \cdot \\
&\quad \left(1 - \psi_i\left(\{x_j\}_{j=1}^{i-1}\right)\right) d\lambda_{X_i|X_1 \dots X_{i-1}} d\lambda_{(X_1, \dots, X_{i-1})} \\
&= \sum_{i=k+1}^{k+r} \int_{\mathbb{R}^{i-1}} \left(1 - \psi_i\left(\{x_j\}_{j=1}^{i-1}\right)\right) \cdot \\
&\quad \int_{\mathbb{R}} (\phi(x_i) - \phi(x_{i-1})) d\lambda_{X_i|X_1 \dots X_{i-1}} d\lambda_{(X_1, \dots, X_{i-1})} \tag{27.5}
\end{aligned}$$

By Jensen's inequality, Lemma 27.3.4 and that this is a sub-martingale, and that ϕ is increasing in addition to being convex,

$$\int_{\mathbb{R}} \phi(x_i) d\lambda_{X_i|X_1 \dots X_{i-1}} \geq \phi\left(\int_{\mathbb{R}} x_i d\lambda_{X_i|X_1 \dots X_{i-1}}\right) \geq \phi(x_{i-1}).$$

Therefore, from 27.5, $E\left(\sum_{i=k+1}^{k+r} (\phi(X_i) - \phi(X_{i-1}))(1 - Y_i)\right) \geq$

$$\sum_{i=k+1}^{k+r} \int_{\mathbb{R}^{i-1}} \left(1 - \psi_i\left(\{x_j\}_{j=1}^{i-1}\right)\right) (\phi(x_{i-1}) - \phi(x_{i-1})) d\lambda_{(X_1, \dots, X_{i-1})} = 0$$

Now let the unbroken strings of ones for $\{Y_i(\omega)\}$ be

$$\{k_1, \dots, k_1 + r_1\}, \{k_2, \dots, k_2 + r_2\}, \dots, \{k_m, \dots, k_m + r_m\} \tag{27.6}$$

where $m = V(\omega) \equiv$ the number of unbroken strings of ones in the sequence $\{Y_i(\omega)\}$. By Corollary 27.3.7 $V(\omega) \geq U_{[a,b]}(\omega)$.

$$\begin{aligned}
&\phi(X_n(\omega)) - \phi(X_1(\omega)) \\
&= \sum_{k=1}^n (\phi(X_k(\omega)) - \phi(X_{k-1}(\omega))) Y_k(\omega) \\
&\quad + \sum_{k=1}^n (\phi(X_k(\omega)) - \phi(X_{k-1}(\omega))) (1 - Y_k(\omega)).
\end{aligned}$$

Summing the first sum over the unbroken strings of ones (the terms in which $Y_i(\omega) = 0$ contribute nothing), and observing that for $x > a$, $\phi(x) = x$,

$$\begin{aligned}
&\phi(X_n(\omega)) - \phi(X_1(\omega)) \geq U_{[a,b]}(\omega)(b - a) + 0 + \\
&\quad \sum_{k=1}^n (\phi(X_k(\omega)) - \phi(X_{k-1}(\omega))) (1 - Y_k(\omega)) \tag{27.7}
\end{aligned}$$

where the zero on the right side results from a string of ones which does not produce an upcrossing. It is here that we use $\phi(x) \geq a$. Such a string begins with $\phi(X_k(\omega)) = a$ and results in an expression of the form $\phi(X_{k+m}(\omega)) - \phi(X_k(\omega)) \geq 0$ since $\phi(X_{k+m}(\omega)) \geq a$.

If we had not replaced X_k with $\phi(X_k)$, it would have been possible for $\phi(X_{k+m}(\omega))$ to be less than a and the zero in the above could have been a negative number.

Therefore from 27.7,

$$\begin{aligned} (b-a)E(U_{[a,b]}) &\leq E(\phi(X_n) - \phi(X_1)) \leq E(\phi(X_n) - a) \\ &= E((X_n - a)^+) \leq |a| + E(|X_n|) \blacksquare \end{aligned}$$

With this estimate, the amazing sub-martingale convergence theorem follows. This incredible theorem says that a bounded in L^1 sub-martingale must converge a.e.

Theorem 27.3.9 (*sub-martingale convergence theorem*) Let $\{X_i\}_{i=1}^\infty$ be a sub-martingale with $K \equiv \sup\{E(|X_n|) : n \geq 1\} < \infty$. Then there exists a random variable X_∞ , such that $E(|X_\infty|) \leq K$ and $\lim_{n \rightarrow \infty} X_n(\omega) = X_\infty(\omega)$ a.e.

Proof: Let $a, b \in \mathbb{Q}$ and let $a < b$. Let $U_{[a,b]}^n(\omega)$ be the number of upcrossings of $\{X_i(\omega)\}_{i=1}^n$. Then let

$$U_{[a,b]}(\omega) \equiv \lim_{n \rightarrow \infty} U_{[a,b]}^n(\omega) = \text{number of upcrossings of } \{X_i\}.$$

By the upcrossing lemma, $E(U_{[a,b]}^n) \leq \frac{E(|X_n|) + |a|}{b-a} \leq \frac{K+|a|}{b-a}$ and so by the monotone convergence theorem, $E(U_{[a,b]}) \leq \frac{K+|a|}{b-a} < \infty$ which shows $U_{[a,b]}(\omega)$ is finite a.e., for all $\omega \notin S_{[a,b]}$ where $P(S_{[a,b]}) = 0$. Define $S \equiv \cup\{S_{[a,b]} : a, b \in \mathbb{Q}, a < b\}$. Then $P(S) = 0$ and if $\omega \notin S$, $\{X_k\}_{k=1}^\infty$ has only finitely many upcrossings of every interval having rational endpoints. Thus, for $\omega \notin S$,

$$\limsup_{k \rightarrow \infty} X_k(\omega) = \liminf_{k \rightarrow \infty} X_k(\omega) = \lim_{k \rightarrow \infty} X_k(\omega) \equiv X_\infty(\omega).$$

Letting $X_\infty(\omega) = 0$ for $\omega \in S$, Fatou's lemma implies

$$\int_{\Omega} |X_\infty| dP = \int_{\Omega} \liminf_{n \rightarrow \infty} |X_n| dP \leq \liminf_{n \rightarrow \infty} \int_{\Omega} |X_n| dP \leq K \blacksquare$$

27.4 Characteristic Functions and Independence

There is a way to tell if random vectors are independent by using their characteristic functions.

Proposition 27.4.1 If X_i is a random vector having values in \mathbb{R}^{p_i} , then the random vectors are independent if and only if

$$E(e^{iP}) = \prod_{j=1}^n E(e^{it_j \cdot X_j})$$

where $P \equiv \sum_{j=1}^n t_j \cdot X_j$ for $t_j \in \mathbb{R}^{p_j}$.

The proof of this proposition will depend on the following lemma.

Lemma 27.4.2 Let \mathbf{Y} be a random vector with values in \mathbb{R}^p and let f be bounded and measurable with respect to the Radon measure $\lambda_{\mathbf{Y}}$, and satisfy

$$\int f(\mathbf{y}) e^{it \cdot \mathbf{y}} d\lambda_{\mathbf{Y}} = 0$$

for all $t \in \mathbb{R}^p$. Then $f(\mathbf{y}) = 0$ for $\lambda_{\mathbf{Y}}$ a.e. \mathbf{y} .

Proof: You could write the following for $\phi \in \mathcal{G}$

$$\int \phi(t) \int f(\mathbf{y}) e^{it \cdot \mathbf{y}} d\lambda_{\mathbf{Y}} dt = 0 = \int f(\mathbf{y}) \left(\int \phi(t) e^{it \cdot \mathbf{y}} dt \right) d\lambda_{\mathbf{Y}}$$

Recall that the inverse Fourier transform maps \mathcal{G} onto \mathcal{G} . Hence $\int f(\mathbf{y}) \psi(\mathbf{y}) d\lambda_{\mathbf{Y}} = 0$ for all $\psi \in \mathcal{G}$. Thus this is also so for every $\psi \in C_0^\infty(\mathbb{R}^p) \supseteq C_c^\infty(\mathbb{R}^p)$ by an obvious application of the Stone Weierstrass theorem. Let $\{\phi_k\}$ be a sequence of functions in $C_c^\infty(\mathbb{R}^p)$ which converges to

$$\text{sgn}(f) \equiv \begin{cases} \bar{f}/|f| & \text{if } f \neq 0 \\ 0 & \text{if } f = 0 \end{cases}$$

pointwise and in $L^1(\mathbb{R}^p, \lambda_{\mathbf{Y}})$, each $|\phi_k| \leq 2$. Then for any $\psi \in C_0^\infty(\mathbb{R}^p)$,

$$0 = \int f(\mathbf{y}) \phi_n(\mathbf{y}) \psi(\mathbf{y}) d\lambda_{\mathbf{Y}} \rightarrow \int |f(\mathbf{y})| \psi(\mathbf{y}) d\lambda_{\mathbf{Y}}$$

Also, the above holds for any $\psi \in C_c(\mathbb{R}^p)$ as can be seen by taking such a ψ and convolving with a mollifier. By the Riesz representation theorem, $f(\mathbf{y}) = 0$ $\lambda_{\mathbf{Y}}$ a.e. (The measure $\mu(E) \equiv \int_E |f(\mathbf{y})| d\lambda_{\mathbf{Y}}$ equals 0.) ■

Proof of the proposition: If the X_j are independent, the formula follows from Lemma 27.2.6 and Lemma 27.2.4.

Now suppose the formula holds. Thus $\prod_{j=1}^n E(e^{it_j \cdot X_j}) =$

$$\begin{aligned} & \int_{\mathbb{R}^{pn}} \cdots \int_{\mathbb{R}^{p2}} \int_{\mathbb{R}^{p1}} e^{it_1 \cdot x_1} e^{it_2 \cdot x_2} \cdots e^{it_n \cdot x_n} d\lambda_{X_1} d\lambda_{X_2} \cdots d\lambda_{X_n} = E(e^{iP}) \\ &= \int_{\mathbb{R}^{pn}} \cdots \int_{\mathbb{R}^{p2}} \int_{\mathbb{R}^{p1}} e^{it_1 \cdot x_1} e^{it_2 \cdot x_2} \cdots e^{it_n \cdot x_n} d\lambda_{X_1|x_2 \cdots x_n} d\lambda_{X_2|x_3 \cdots x_n} \cdots d\lambda_{X_n}. \end{aligned} \quad (27.8)$$

Then from the above Lemma 27.4.2, the following equals 0 for λ_{X_n} a.e. x_n .

$$\begin{aligned} & \int_{\mathbb{R}^{p(n-1)}} \cdots \int_{\mathbb{R}^{p2}} \int_{\mathbb{R}^{p1}} e^{it_1 \cdot x_1} e^{it_2 \cdot x_2} \cdots e^{it_{n-1} \cdot x_{n-1}} d\lambda_{X_1} d\lambda_{X_2} \cdots d\lambda_{X_{n-1}} - \\ & \int_{\mathbb{R}^{p(n-1)}} \cdots \int_{\mathbb{R}^{p2}} \int_{\mathbb{R}^{p1}} e^{it_1 \cdot x_1} e^{it_2 \cdot x_2} \\ & \cdots e^{it_{n-1} \cdot x_{n-1}} d\lambda_{X_1|x_2 \cdots x_n} d\lambda_{X_2|x_3 \cdots x_n} \cdots d\lambda_{X_{n-1}|x_n} \end{aligned}$$

Let $t_i = 0$ for $i = 1, 2, \dots, n-2$. Then this implies

$$\int_{\mathbb{R}^{p(n-1)}} e^{it_{n-1} \cdot x_{n-1}} d\lambda_{X_{n-1}} = \int_{\mathbb{R}^{p(n-1)}} e^{it_{n-1} \cdot x_{n-1}} d\lambda_{X_{n-1}|x_n}$$

By the fact that the characteristic function determines the distribution measure, Theorem 27.1.4, it follows that for these x_n off a set of λ_{X_n} measure zero, $\lambda_{X_{n-1}} = \lambda_{X_{n-1}|x_n}$. Returning to 27.8, one can replace $\lambda_{X_{n-1}|x_n}$ with $\lambda_{X_{n-1}}$ to obtain

$$\begin{aligned} & \int_{\mathbb{R}^{p_n}} \cdots \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_1}} e^{it_1 \cdot x_1} e^{it_2 \cdot x_2} \cdots e^{it_n \cdot x_n} d\lambda_{X_1} d\lambda_{X_2} \cdots d\lambda_{X_{n-1}} d\lambda_{X_n} \\ &= \int_{\mathbb{R}^{p_n}} \cdots \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_1}} e^{it_1 \cdot x_1} e^{it_2 \cdot x_2} \cdots e^{it_n \cdot x_n} \cdot \\ & \quad d\lambda_{X_1|x_2 \dots x_n} d\lambda_{X_2|x_3 \dots x_n} \cdots d\lambda_{X_{n-1}} d\lambda_{X_n} \end{aligned}$$

Next let $t_n = 0$ and applying the above Lemma 27.4.2 again, this implies that for $\lambda_{X_{n-1}}$ a.e. x_{n-1} , the following equals 0.

$$\begin{aligned} & \int_{\mathbb{R}^{p_{n-2}}} \cdots \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_1}} e^{it_1 \cdot x_1} e^{it_2 \cdot x_2} \cdots e^{it_{n-2} \cdot x_{n-2}} d\lambda_{X_1} d\lambda_{X_2} \cdots d\lambda_{X_{n-2}} - \\ & \int_{\mathbb{R}^{p_{n-2}}} \cdots \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_1}} e^{it_1 \cdot x_1} e^{it_2 \cdot x_2} \cdots e^{it_{n-2} \cdot x_{n-2}} \cdot \\ & \quad d\lambda_{X_1|x_2 \dots x_n} d\lambda_{X_2|x_3 \dots x_n} \cdots d\lambda_{X_{n-2}|x_n x_{n-1}} \end{aligned}$$

Let $t_i = 0$ for $i = 1, 2, \dots, n-3$. Then you obtain

$$\int_{\mathbb{R}^{p_{n-2}}} e^{it_{n-2} \cdot x_{n-2}} d\lambda_{X_{n-2}} = \int_{\mathbb{R}^{p_{n-2}}} e^{it_{n-2} \cdot x_{n-2}} d\lambda_{X_{n-2}|x_n x_{n-1}}$$

and so $\lambda_{X_{n-2}} = \lambda_{X_{n-2}|x_n x_{n-1}}$ for x_{n-1} off a set of $\lambda_{X_{n-1}}$ measure zero. Continuing this way, it follows that

$$\lambda_{X_{n-k}} = \lambda_{X_{n-k}|x_n x_{n-1} \cdots x_{n-k+1}}$$

for x_{n-k+1} off a set of $\lambda_{X_{n-k+1}}$ measure zero. Thus if E is Borel in $\mathbb{R}^{p_{n-1}} \times \cdots \times \mathbb{R}^{p_1}$,

$$\begin{aligned} & \int_{\mathbb{R}^{p_n} \times \cdots \times \mathbb{R}^{p_1}} \mathcal{X}_E d\lambda_{(X_1 \cdots X_n)} = \\ & \int_{\mathbb{R}^{p_n}} \cdots \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_1}} \mathcal{X}_E d\lambda_{X_1|x_2 \dots x_n} d\lambda_{X_2|x_3 \dots x_n} \cdots d\lambda_{X_{n-1}|x_n} d\lambda_{X_n} \\ & \int_{\mathbb{R}^{p_n}} \cdots \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_1}} \mathcal{X}_E d\lambda_{X_1|x_2 \dots x_n} d\lambda_{X_2|x_3 \dots x_n} \cdots d\lambda_{X_{n-1}} d\lambda_{X_n} \\ & \quad \vdots \\ &= \int_{\mathbb{R}^{p_n}} \cdots \int_{\mathbb{R}^{p_2}} \int_{\mathbb{R}^{p_1}} \mathcal{X}_E d\lambda_{X_1} d\lambda_{X_2} \cdots d\lambda_{X_n} \end{aligned}$$

One could achieve this iterated integral in any order by similar arguments to the above. By Definition 27.2.2 and the discussion which follows, this implies that the random variables X_i are independent. ■

Here is another proof of the Doob Dynkin lemma based on differentiation theory.

Lemma 27.4.3 Suppose $\mathbf{X}, \mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_k$ are random vectors \mathbf{X} having values in \mathbb{R}^n and \mathbf{Y}_j having values in \mathbb{R}^{p_j} and

$$\mathbf{X}, \mathbf{Y}_j \in L^1(\Omega).$$

Suppose \mathbf{X} is $\sigma(\mathbf{Y}_1, \dots, \mathbf{Y}_k)$ measurable. Thus

$$\{\mathbf{X}^{-1}(E) : E \text{ Borel}\} \subseteq \left\{ (\mathbf{Y}_1, \dots, \mathbf{Y}_k)^{-1}(F) : F \text{ is Borel in } \prod_{j=1}^k \mathbb{R}^{p_j} \right\}$$

Then there exists a Borel function, $\mathbf{g} : \prod_{j=1}^k \mathbb{R}^{p_j} \rightarrow \mathbb{R}^n$ such that

$$\mathbf{X} = \mathbf{g}(\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_k).$$

Proof: For the sake of brevity, denote by \mathbf{Y} the vector $(\mathbf{Y}_1, \dots, \mathbf{Y}_k)$ and by \mathbf{y} the vector $(\mathbf{y}_1, \dots, \mathbf{y}_k)$ and let $\prod_{j=1}^k \mathbb{R}^{p_j} \equiv \mathbb{R}^P$. For E a Borel set of \mathbb{R}^n ,

$$\begin{aligned} \int_{\mathbf{Y}^{-1}(E)} \mathbf{X} dP &= \int_{\mathbb{R}^n \times \mathbb{R}^P} \mathcal{X}_{\mathbb{R}^n \times E}(\mathbf{x}, \mathbf{y}) \mathbf{x} d\lambda_{(\mathbf{X}, \mathbf{Y})} \\ &= \int_E \int_{\mathbb{R}^n} \mathbf{x} d\lambda_{\mathbf{X}|\mathbf{y}} d\lambda_{\mathbf{Y}}. \end{aligned} \quad (27.9)$$

Consider the function $\mathbf{y} \rightarrow \int_{\mathbb{R}^n} \mathbf{x} d\lambda_{\mathbf{X}|\mathbf{y}}$. Since $d\lambda_{\mathbf{Y}}$ is a Radon measure having inner and outer regularity, it follows the above function is equal to a Borel function for $\lambda_{\mathbf{Y}}$ a.e. \mathbf{y} . This function will be denoted by \mathbf{g} . Then from 27.9

$$\begin{aligned} \int_{\mathbf{Y}^{-1}(E)} \mathbf{X} dP &= \int_E \mathbf{g}(\mathbf{y}) d\lambda_{\mathbf{Y}} = \int_{\mathbb{R}^P} \mathcal{X}_E(\mathbf{y}) \mathbf{g}(\mathbf{y}) d\lambda_{\mathbf{Y}} \\ &= \int_{\Omega} \mathcal{X}_E(\mathbf{Y}(\omega)) \mathbf{g}(\mathbf{Y}(\omega)) dP = \int_{\mathbf{Y}^{-1}(E)} \mathbf{g}(\mathbf{Y}(\omega)) dP \end{aligned}$$

and since $\mathbf{Y}^{-1}(E)$ is an arbitrary element of $\sigma(\mathbf{Y})$, this shows that since \mathbf{X} is $\sigma(\mathbf{Y})$ measurable, $\mathbf{X} = \mathbf{g}(\mathbf{Y})$ P a.e. ■

What about the case where \mathbf{X} is not necessarily measurable in $\sigma(\mathbf{Y}_1, \dots, \mathbf{Y}_k)$?

Lemma 27.4.4 There exists a unique function $\mathbf{Z}(\omega)$ which satisfies

$$\int_F \mathbf{X}(\omega) dP = \int_F \mathbf{Z}(\omega) dP$$

for all

$$F \in \sigma(\mathbf{Y}_1, \dots, \mathbf{Y}_k)$$

such that \mathbf{Z} is $\sigma(\mathbf{Y}_1, \dots, \mathbf{Y}_k)$ measurable. It is denoted by

$$E(\mathbf{X} | \sigma(\mathbf{Y}_1, \dots, \mathbf{Y}_k))$$

Proof: It is like the above. Letting E be a Borel set in \mathbb{R}^P ,

$$\int_{\mathbf{Y}^{-1}(E)} \mathbf{X} dP = \int_{\mathbb{R}^n \times \mathbb{R}^P} \mathcal{X}_{\mathbb{R}^n \times E}(\mathbf{x}, \mathbf{y}) \mathbf{x} d\lambda_{(\mathbf{X}, \mathbf{Y})} = \int_E \int_{\mathbb{R}^n} \mathbf{x} d\lambda_{\mathbf{X}|\mathbf{y}} d\lambda_{\mathbf{Y}}.$$

Now let $g(y) \equiv E(X|y_1, \dots, y_k)$ be a Borel representative of $\int_{\mathbb{R}^n} x d\lambda_{X|y}$. It follows $\omega \rightarrow g(Y(\omega)) = E(X|Y_1(\omega), \dots, Y_k(\omega))$ is $\sigma(Y_1, \dots, Y_k)$ measurable because by definition $\omega \rightarrow Y(\omega)$ is $\sigma(Y_1, \dots, Y_k)$ measurable and a Borel measurable function composed with a measurable one is still measurable. It follows that for all E Borel in \mathbb{R}^p ,

$$\begin{aligned} \int_{Y^{-1}(E)} X dP &= \int_E E(X|y_1, \dots, y_k) d\lambda_Y \\ &= \int_{Y^{-1}(E)} E(X|Y_1(\omega), \dots, Y_k(\omega)) dP \end{aligned}$$

so $Z(\omega) = E(X|Y_1(\omega), \dots, Y_k(\omega))$ works because a generic set of $\sigma(Y_1, \dots, Y_k)$ is $Y^{-1}(E)$ for E a Borel set in \mathbb{R}^p . If both Z, Z_1 work, then for all

$$F \in \sigma(Y_1, \dots, Y_k),$$

$$\int_F (Z - Z_1) dP = 0$$

Since F is arbitrary, some routine computations show $Z = Z_1$ a.e. ■

Observation 27.4.5 *Note that a.e.*

$$E(X|Y_1(\omega), \dots, Y_k(\omega)) = E(X|\sigma(Y_1, \dots, Y_k))$$

where the one on the left is the expected value of X given values of $Y_j(\omega)$. This one corresponds to the sort of thing we say in words. The one on the right is an abstract concept which is usually obtained using the Radon Nikodym theorem and its description is given in the lemma. This lemma shows that its meaning is really to take the expected value of X given values for the Y_k .

27.5 Characteristic Functions for Measures

Recall the characteristic function for a random variable having values in \mathbb{R}^p . I will give a review of this to begin with. Then the concept will be generalized to random variables (vectors) which have values in a real separable Banach space.

Definition 27.5.1 *Let X be a random variable. The characteristic function is*

$$\phi_X(t) \equiv E(e^{it \cdot X}) \equiv \int_{\Omega} e^{it \cdot X(\omega)} dP = \int_{\mathbb{R}^p} e^{it \cdot x} d\lambda_X$$

the last equation holding by Proposition 26.1.12 on Page 717.

Recall the following fundamental lemma and definition, Lemma 13.2.4 on Page 379.

Definition 27.5.2 *For $T \in \mathcal{G}^*$, define $FT, F^{-1}T \in \mathcal{G}^*$ by*

$$FT(\phi) \equiv T(F\phi), F^{-1}T(\phi) \equiv T(F^{-1}\phi)$$

Lemma 27.5.3 *F and F^{-1} are both one to one, onto, and are inverses of each other.*

The main result on characteristic functions is the following in Theorem 27.1.4 on Page 735 which is stated here for convenience.

Theorem 27.5.4 Let X and Y be random vectors with values in \mathbb{R}^p and suppose $E(e^{it \cdot X}) = E(e^{it \cdot Y})$ for all $t \in \mathbb{R}^p$. Then $\lambda_X = \lambda_Y$.

I want to do something similar for random variables which have values in a separable real Banach space, E instead of \mathbb{R}^p .

Corollary 27.5.5 Let \mathcal{K} be a π system of subsets of Ω and suppose two probability measures, μ and ν defined on $\sigma(\mathcal{K})$ are equal on \mathcal{K} . Then $\mu = \nu$.

Proof: This follows from the Lemma 9.3.2 on Page 243. Let

$$\mathcal{G} \equiv \{E \in \sigma(\mathcal{K}) : \mu(E) = \nu(E)\}$$

Then $\mathcal{K} \subseteq \mathcal{G}$, since μ and ν are both probability measures, it follows that if $E \in \mathcal{G}$, then so is E^C . Since these are measures, if $\{A_i\}$ is a sequence of disjoint sets from \mathcal{G} then

$$\mu(\cup_{i=1}^{\infty} A_i) = \sum_i \mu(A_i) = \sum_i \nu(A_i) = \nu(\cup_{i=1}^{\infty} A_i)$$

and so from Lemma 9.3.2, $\mathcal{G} = \sigma(\mathcal{K})$. ■

Next recall the following fundamental lemma used to prove Pettis' theorem. It is proved on Page 649 but is stated here for convenience.

Lemma 27.5.6 If E is a separable Banach space with B' the closed unit ball in E' , then there exists a sequence $\{f_n\}_{n=1}^{\infty} \equiv D' \subseteq B'$ with the property that for every $x \in E$,

$$\|x\| = \sup_{f \in D'} |f(x)|$$

Definition 27.5.7 Let E be a separable real Banach space. A cylindrical set is one which is of the form

$$\{x \in E : x_i^*(x) \in \Gamma_i, i = 1, 2, \dots, m\}$$

where here $x_i^* \in E'$ and Γ_i is a Borel set in \mathbb{R} .

It is obvious that \emptyset is a cylindrical set and that the intersection of two cylindrical sets is another cylindrical set. Thus the cylindrical sets form a π system. What is the smallest σ algebra containing the cylindrical sets? It is the Borel sets of E . This is a special case of Lemma 26.5.2. Recall why this was. Letting $\{f_n\}_{n=1}^{\infty} = D'$ be the sequence of Lemma 27.5.6 it follows that

$$\begin{aligned} & \{x \in E : \|x - a\| \leq \delta\} \\ &= \left\{ x \in E : \sup_{f \in D'} |f(x - a)| \leq \delta \right\} = \left\{ x \in E : \sup_{f \in D'} |f(x) - f(a)| \leq \delta \right\} \\ &= \cap_{n=1}^{\infty} \left\{ x \in E : f_n(x) \in \overline{B(f_n(a), \delta)} \right\} \end{aligned}$$

which yields a countable intersection of cylindrical sets. It follows the smallest σ algebra containing the cylindrical sets contains the closed balls and hence the open balls and consequently the open sets and so it contains the Borel sets. However, each cylindrical set is a Borel set and so in fact this σ algebra equals $\mathcal{B}(E)$.

From Corollary 27.5.5 it follows that two probability measures which are equal on the cylindrical sets are equal on the Borel sets $\mathcal{B}(E)$.

Definition 27.5.8 Let μ be a probability measure on $\mathcal{B}(E)$ where E is a real separable Banach space. Then for $x^* \in E'$,

$$\phi_\mu(x^*) \equiv \int_E e^{ix^*(x)} d\mu(x).$$

ϕ_μ is called the characteristic function for the measure μ .

Note this is a little different than earlier when the symbol $\phi_X(t)$ was used and X was a random variable. Here the focus is more on the measure than a random variable X such that its distribution measure is μ . It might appear this is a more general concept but in fact this is not the case. You could just consider the separable Banach space or Polish space with the Borel σ algebra as your probability space and then consider the identity map as a random variable having the given measure as a distribution measure. Of course a major result is the one which says that the characteristic function determines the measures.

Theorem 27.5.9 Let μ and ν be two probability measures on $\mathcal{B}(E)$ where E is a separable real Banach space. Suppose

$$\phi_\mu(x^*) = \phi_\nu(x^*)$$

for all $x^* \in E'$. Then $\mu = \nu$.

Proof: It suffices to verify that $\mu(A) = \nu(A)$ for all $A \in \mathcal{K}$ where \mathcal{K} is the set of cylindrical sets. Fix $g_n \in (E')^n$. Thus the two measures are equal if for all such g_n , $n \in \mathbb{N}$,

$$\mu(g_n^{-1}(B)) = \nu(g_n^{-1}(B))$$

for B a Borel set in \mathbb{R}^n . Of course, for such a choice of $g_n \in (E')^n$, there are measures defined on the Borel sets of \mathbb{R}^n μ_n and ν_n which are given by

$$\mu_n(B) \equiv \mu(g_n^{-1}(B)), \quad \nu_n(B) \equiv \nu(g_n^{-1}(B))$$

and so it suffices to verify that these two measures are equal. So what are their characteristic functions? Note that g_n is a random variable taking E to \mathbb{R}^n and μ_n, ν_n are just the probability distribution measures of this random variable. Therefore,

$$\phi_{\mu_n}(t) \equiv \int_{\mathbb{R}^n} e^{it \cdot s} d\mu_n = \int_E e^{it \cdot g_n(x)} d\mu$$

Similarly,

$$\phi_{\nu_n}(t) \equiv \int_{\mathbb{R}^n} e^{it \cdot s} d\nu_n = \int_E e^{it \cdot g_n(x)} d\nu$$

Now $t \cdot g_n \in E'$ and so by assumption, the two ends of the above are equal. Hence $\phi_{\mu_n}(t) = \phi_{\nu_n}(t)$ and so by Theorem 27.1.6, $\mu_n = \nu_n$ which, as shown above, implies $\mu = \nu$. ■

27.6 Independence in Banach Space

I will consider the relation between the characteristic function and independence of random variables having values in a Banach space. Recall an earlier proposition which relates independence of random vectors with characteristic functions. It is proved starting on Page 744.

Proposition 27.6.1 *Let $\{X_k\}_{k=1}^n$ be random vectors such that X_k has values in \mathbb{R}^{p_k} . Then the random vectors are independent if and only if*

$$E(e^{iP}) = \prod_{j=1}^n E(e^{it_j \cdot X_j})$$

where $P \equiv \sum_{j=1}^n t_j \cdot X_j$ for $t_j \in \mathbb{R}^{p_j}$.

It turns out there is a generalization of the above proposition to the case where the random variables have values in a real separable Banach space. Before proving this recall an earlier theorem which had to do with reducing to the case where the random variables had values in \mathbb{R}^n , Theorem 26.6.1. It is restated here for convenience.

Theorem 27.6.2 *The random variables $\{X_i\}_{i \in I}$ are independent if whenever*

$$\{i_1, \dots, i_n\} \subseteq I,$$

m_{i_1}, \dots, m_{i_n} are positive integers, and $g_{m_{i_1}}, \dots, g_{m_{i_n}}$ are in

$$(E')^{m_{i_1}}, \dots, (E')^{m_{i_n}}$$

respectively, $\{g_{m_{i_j}} \circ X_{i_j}\}_{j=1}^n$ are independent random vectors having values in

$$\mathbb{R}^{m_{i_1}}, \dots, \mathbb{R}^{m_{i_n}}$$

respectively.

Now here is the theorem about independence and the characteristic functions.

Theorem 27.6.3 *Let $\{X_k\}_{k=1}^n$ be random variables such that X_k has values in E_k , a real separable Banach space. Then the random variables are independent if and only if*

$$E(e^{iP}) = \prod_{j=1}^n E(e^{it_j^*(X_j)})$$

where $P \equiv \sum_{j=1}^n t_j^*(X_j)$ for $t_j^* \in E'_j$.

Proof: If the random variables are independent, then so are the random variables, $t_j^*(X_j)$ and so the equation follows.

The interesting case is when the equation holds.

It suffices to consider only the case where each $E_k = E$. This is because you can consider each X_j to have values in $\prod_{k=1}^n E_k$ by letting X_j take its values in the j^{th} component of the product and 0 in the other components. Can you draw the conclusion the random variables are independent? By Theorem 26.6.1, it suffices to show the random variables $\{g_{m_k} \circ X_k\}_{k=1}^n$ are independent where $g_{m_k} = (x_1^*, \dots, x_{m_k}^*) \in (E')^{m_k}$. This happens if whenever $t_{m_k} \in \mathbb{R}^{m_k}$ and

$$P = \sum_{k=1}^n t_{m_k} \cdot (g_{m_k} \circ X_k),$$

it follows

$$E(e^{iP}) = \prod_{k=1}^n E\left(e^{it_{m_k} \cdot (g_{m_k} \circ X_k)}\right). \quad (27.10)$$

However, the expression on the right in 27.10 equals

$$\prod_{k=1}^n E\left(e^{i(t_{m_k} \cdot g_{m_k}) \circ X_k}\right)$$

and $t_{m_k} \cdot g_{m_k} \equiv \sum_{j=1}^{m_k} t_j x_j^* \in E'$. Also the expression on the left equals

$$E\left(e^{i \sum_{k=1}^n t_{m_k} \cdot g_{m_k} \circ X_k}\right)$$

Therefore, by assumption, 27.10 holds. ■

There is an obvious corollary which is useful.

Corollary 27.6.4 *Let $\{X_k\}_{k=1}^n$ be random variables such that X_k has values in E_k , a real separable Banach space. Then the random variables are independent if and only if*

$$E(e^{iP}) = \prod_{j=1}^n E(e^{it_j^*(X_j)})$$

where $P \equiv \sum_{j=1}^n t_j^*(X_j)$ for $t_j^* \in M_j$ where M_j is a dense subset of E'_j .

Proof: The easy direction follows from Theorem 27.6.3. Suppose then the above equation holds for all $t_j^* \in M_j$. Then let $t_j^* \in E'$ and let $\{t_{n_j}^*\}$ be a sequence in M_j such that $\lim_{n \rightarrow \infty} t_{n_j}^* = t_j^*$ in E' . Then define

$$P \equiv \sum_{j=1}^n t_j^* X_j, \quad P_n \equiv \sum_{j=1}^n t_{n_j}^* X_j.$$

It follows

$$E(e^{iP}) = \lim_{n \rightarrow \infty} E(e^{iP_n}) = \lim_{n \rightarrow \infty} \prod_{j=1}^n E(e^{it_{n_j}^*(X_j)}) = \prod_{j=1}^n E(e^{it_j^*(X_j)}) \quad \blacksquare$$

27.7 Convolution and Sums

Lemma 26.1.9 on Page 717 makes possible a definition of convolution of two probability measures defined on $\mathcal{B}(E)$ where E is a separable Banach space. I will first show a little theorem about density of continuous functions in $L^p(E)$ and then define the convolution of two finite measures. First here is a simple technical lemma.

Lemma 27.7.1 *Suppose K is a compact subset of U an open set in E a metric space. Then there exists $\delta > 0$ such that*

$$\text{dist}(x, K) + \text{dist}(x, U^C) \geq \delta \text{ for all } x \in E.$$

Proof: For each $x \in K$, there exists a ball, $B(x, \delta_x)$ such that $B(x, 3\delta_x) \subseteq U$. Finitely many of these balls cover K because K is compact, say $\{B(x_i, \delta_{x_i})\}_{i=1}^m$. Let

$$0 < \delta < \min(\delta_{x_i} : i = 1, 2, \dots, m).$$

Now pick any $x \in K$. Then $x \in B(x_i, \delta_{x_i})$ for some x_i and so

$$B(x, \delta) \subseteq B(x_i, 2\delta_{x_i}) \subseteq U.$$

Therefore, for any $x \in K$, $\text{dist}(x, U^C) \geq \delta$. If $x \in B(x_i, 2\delta_{x_i})$ for some x_i , it follows that $\text{dist}(x, U^C) \geq \delta$ because then $B(x, \delta) \subseteq B(x_i, 3\delta_{x_i}) \subseteq U$. If $x \notin B(x_i, 2\delta_{x_i})$ for any of the x_i , then $x \notin B(y, \delta)$ for any $y \in K$ because all these sets are contained in some $B(x_i, 2\delta_{x_i})$. Consequently $\text{dist}(x, K) \geq \delta$. ■

From this lemma, there is an easy corollary.

Corollary 27.7.2 *Suppose K is a compact subset of U , an open set in E a metric space. Then there exists a uniformly continuous function f defined on all of E , having values in $[0, 1]$ such that $f(x) = 0$ if $x \notin U$ and $f(x) = 1$ if $x \in K$.*

Proof: Consider

$$f(x) \equiv \frac{\text{dist}(x, U^C)}{\text{dist}(x, U^C) + \text{dist}(x, K)}.$$

Then some algebra yields

$$|f(x) - f(x')| \leq \frac{1}{\delta} (|\text{dist}(x, U^C) - \text{dist}(x', U^C)| + |\text{dist}(x, K) - \text{dist}(x', K)|)$$

where δ is the constant of Lemma 27.7.1. Now it is a general fact that

$$|\text{dist}(x, S) - \text{dist}(x', S)| \leq d(x, x').$$

See Proposition 3.6.6. Therefore, $|f(x) - f(x')| \leq \frac{2}{\delta} d(x, x')$ and this proves the corollary. ■

Now suppose μ is a finite measure defined on the Borel sets of a separable Banach space E . It was shown above that μ is inner and outer regular. Lemma 26.1.9 on Page 717 shows that μ is inner regular in the usual sense with respect to compact sets. This makes possible the following theorem.

Theorem 27.7.3 *Let μ be a finite measure on $\mathcal{B}(E)$ where E is a separable Banach space and let $f \in L^p(E; \mu)$. Then for any $\varepsilon > 0$, there exists a uniformly continuous, bounded g defined on E such that*

$$\|f - g\|_{L^p(E)} < \varepsilon.$$

Proof: As usual in such situations, it suffices to consider only $f \geq 0$. Then by Theorem 9.1.6 on Page 239 and an application of the monotone convergence theorem, there exists a simple measurable function,

$$s(x) \equiv \sum_{k=1}^m c_k \chi_{A_k}(x)$$

such that $\|f - s\|_{L^p(E)} < \varepsilon/2$. Now by regularity of μ there exist compact sets, K_k and open sets, V_k such that $2 \sum_{k=1}^m |c_k| \mu(V_k \setminus K)^{1/p} < \varepsilon/2$ and by Corollary 27.7.2 there exist

uniformly continuous functions g_k having values in $[0, 1]$ such that $g_k = 1$ on K_k and 0 on V_k^C . Then consider

$$g(x) = \sum_{k=1}^m c_k g_k(x).$$

This function is bounded and uniformly continuous. Furthermore,

$$\begin{aligned} \|s - g\|_{L^p(E)} &\leq \left(\int_E \left| \sum_{k=1}^m c_k \mathcal{X}_{A_k}(x) - \sum_{k=1}^m c_k g_k(x) \right|^p d\mu \right)^{1/p} \\ &\leq \left(\int_E \left(\sum_{k=1}^m |c_k| |\mathcal{X}_{A_k}(x) - g_k(x)| \right)^p d\mu \right)^{1/p} \leq \sum_{k=1}^m |c_k| \left(\int_E |\mathcal{X}_{A_k}(x) - g_k(x)|^p d\mu \right)^{1/p} \\ &\leq \sum_{k=1}^m |c_k| \left(\int_{V_k \setminus K_k} 2^p d\mu \right)^{1/p} = 2 \sum_{k=1}^m |c_k| \mu(V_k \setminus K_k)^{1/p} < \varepsilon/2. \end{aligned}$$

Therefore,

$$\|f - g\|_{L^p} \leq \|f - s\|_{L^p} + \|s - g\|_{L^p} < \varepsilon/2 + \varepsilon/2 \quad \blacksquare$$

Lemma 27.7.4 *Let $A \in \mathcal{B}(E)$ where μ is a finite measure on $\mathcal{B}(E)$ for E a separable Banach space. Also let $x_i \in E$ for $i = 1, 2, \dots, m$. Then for $\mathbf{x} \in E^m$,*

$$\mathbf{x} \rightarrow \mu \left(A + \sum_{i=1}^m x_i \right), \quad \mathbf{x} \rightarrow \mu \left(A - \sum_{i=1}^m x_i \right)$$

are Borel measurable functions. Furthermore, the above functions are

$$\mathcal{B}(E) \times \cdots \times \mathcal{B}(E)$$

measurable where the above denotes the product measurable sets as described in Theorem 10.14.9 on Page 306.

Proof: First consider the case where $A = U$, an open set. Let

$$\mathbf{y} \in \left\{ \mathbf{x} \in E^m : \mu \left(U + \sum_{i=1}^m x_i \right) > \alpha \right\} \quad (27.11)$$

Then from Lemma 26.1.9 on Page 717 there exists a compact set, $K \subseteq U + \sum_{i=1}^m y_i$ such that $\mu(K) > \alpha$. Then if \mathbf{y}' is close enough to \mathbf{y} , it follows $K \subseteq U + \sum_{i=1}^m y'_i$ also. Therefore, for all \mathbf{y}' close enough to \mathbf{y} ,

$$\mu \left(U + \sum_{i=1}^m y'_i \right) \geq \mu(K) > \alpha.$$

In other words the set described in 27.11 is an open set and so $\mathbf{y} \rightarrow \mu(U + \sum_{i=1}^m y_i)$ is Borel measurable whenever U is an open set in E .

Define a π system, \mathcal{K} to consist of all open sets in E . Then define \mathcal{G} as

$$\left\{ A \in \sigma(\mathcal{K}) = \mathcal{B}(E) : \mathbf{y} \rightarrow \mu \left(A + \sum_{i=1}^m y_i \right) \text{ is Borel measurable} \right\}$$

I just showed $\mathcal{G} \supseteq \mathcal{K}$. Now suppose $A \in \mathcal{G}$. Then

$$\mu \left(A^C + \sum_{i=1}^m y_i \right) = \mu(E) - \mu \left(A + \sum_{i=1}^m y_i \right)$$

and so $A^C \in \mathcal{G}$ whenever $A \in \mathcal{G}$. Next suppose $\{A_i\}$ is a sequence of disjoint sets of \mathcal{G} . Then

$$\mu \left(\left(\bigcup_{i=1}^{\infty} A_i \right) + \sum_{j=1}^m y_j \right) = \mu \left(\bigcup_{i=1}^{\infty} \left(A_i + \sum_{j=1}^m y_j \right) \right) = \sum_{i=1}^{\infty} \mu \left(A_i + \sum_{j=1}^m y_j \right)$$

and so $\bigcup_{i=1}^{\infty} A_i \in \mathcal{G}$ because the above is the sum of Borel measurable functions. By the lemma on π systems, Lemma 9.3.2 on Page 243, it follows $\mathcal{G} = \sigma(\mathcal{K}) = \mathcal{B}(E)$. Similarly, $x \rightarrow \mu \left(A - \sum_{j=1}^m x_j \right)$ is also Borel measurable whenever $A \in \mathcal{B}(E)$. Finally note that

$$\mathcal{B}(E) \times \cdots \times \mathcal{B}(E)$$

contains the open sets of E^m because the separability of E implies the existence of a countable basis for the topology of E^m consisting of sets of the form $\prod_{i=1}^m U_i$ where the U_i come from a countable basis for E . Since every open set is the countable union of sets like the above, each being a measurable box, the open sets are contained in $\mathcal{B}(E) \times \cdots \times \mathcal{B}(E)$ which implies $\mathcal{B}(E^m) \subseteq \mathcal{B}(E) \times \cdots \times \mathcal{B}(E)$ also. ■

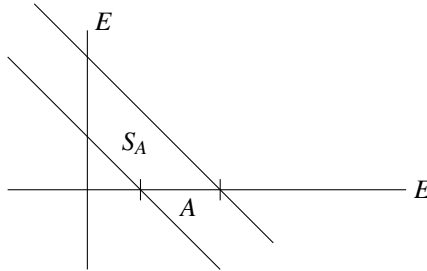
With this lemma, it is possible to define the convolution of two finite measures.

Definition 27.7.5 Let μ and ν be two finite measures on $\mathcal{B}(E)$, for E a separable Banach space. Then define a new measure, $\mu * \nu$ on $\mathcal{B}(E)$ as follows

$$\mu * \nu(A) \equiv \int_E \nu(A - x) d\mu(x).$$

This is well defined because of Lemma 27.7.4 which says that $x \rightarrow \nu(A - x)$ is Borel measurable.

Here is an interesting theorem about convolutions. However, first here is a little lemma. The following picture is descriptive of the set described in the following lemma.



Lemma 27.7.6 For A a Borel set in E , a separable Banach space, define

$$S_A \equiv \{(x, y) \in E \times E : x + y \in A\}$$

Then $S_A \in \mathcal{B}(E) \times \mathcal{B}(E)$, the σ algebra of product measurable sets, the smallest σ algebra which contains all the sets of the form $A \times B$ where A and B are Borel.

Proof: Let \mathcal{K} denote the open sets in E . Then \mathcal{K} is a π system. Let

$$\mathcal{G} \equiv \{A \in \sigma(\mathcal{K}) = \mathcal{B}(E) : S_A \in \mathcal{B}(E) \times \mathcal{B}(E)\}.$$

Then $\mathcal{K} \subseteq \mathcal{G}$ because if $U \in \mathcal{K}$ then S_U is an open set in $E \times E$ and all open sets are in $\mathcal{B}(E) \times \mathcal{B}(E)$ because a countable basis for the topology of $E \times E$ are sets of the form $B \times C$ where B and C come from a countable basis for E . Therefore, $\mathcal{K} \subseteq \mathcal{G}$. Now let $A \in \mathcal{G}$. For $(x, y) \in E \times E$, either $x + y \in A$ or $x + y \notin A$. Hence $E \times E = S_A \cup S_{A^C}$ which shows that if $A \in \mathcal{G}$ then so is A^C . Finally if $\{A_i\}$ is a sequence of disjoint sets of \mathcal{G}

$$S_{\cup_{i=1}^{\infty} A_i} = \cup_{i=1}^{\infty} S_{A_i}$$

and this shows that \mathcal{G} is also closed with respect to countable unions of disjoint sets. Therefore, by the lemma on π systems, Lemma 9.3.2 on Page 243 it follows $\mathcal{G} = \sigma(\mathcal{K}) = \mathcal{B}(E)$. This proves the lemma.

Theorem 27.7.7 *Let μ , ν , and λ be finite measures on $\mathcal{B}(E)$ for E a separable Banach space. Then*

$$\mu * \nu = \nu * \mu \tag{27.12}$$

$$(\mu * \nu) * \lambda = \mu * (\nu * \lambda) \tag{27.13}$$

*If μ is the distribution for an E valued random variable, X and if ν is the distribution for an E valued random variable, Y , and X and Y are independent, then $\mu * \nu$ is the distribution for the random variable, $X + Y$. Also the characteristic function of a convolution equals the product of the characteristic functions.*

Proof: First consider 27.12. Letting $A \in \mathcal{B}(E)$, the following computation holds from Fubini's theorem and Lemma 27.7.6

$$\begin{aligned} \mu * \nu(A) &\equiv \int_E \nu(A - x) d\mu(x) = \int_E \int_E \mathcal{X}_{S_A}(x, y) d\nu(y) d\mu(x) \\ &= \int_E \int_E \mathcal{X}_{S_A}(x, y) d\mu(x) d\nu(y) = \nu * \mu(A). \end{aligned}$$

Next consider 27.13. Using 27.12 whenever convenient,

$$\begin{aligned} (\mu * \nu) * \lambda(A) &\equiv \int_E (\mu * \nu)(A - x) d\lambda(x) \\ &= \int_E \int_E \nu(A - x - y) d\mu(y) d\lambda(x) \end{aligned}$$

while

$$\begin{aligned} \mu * (\nu * \lambda)(A) &\equiv \int_E (\nu * \lambda)(A - y) d\mu(y) \\ &= \int_E \int_E \nu(A - y - x) d\lambda(x) d\mu(y) \\ &= \int_E \int_E \nu(A - y - x) d\mu(y) d\lambda(x). \end{aligned}$$

The necessary product measurability comes from Lemma 27.7.4.

Recall

$$(\mu * \nu)(A) \equiv \int_E \nu(A-x) d\mu(x).$$

Therefore, if s is a simple function, $s(x) = \sum_{k=1}^n c_k \mathcal{X}_{A_k}(x)$,

$$\begin{aligned} \int_E s d(\mu * \nu) &= \sum_{k=1}^n c_k \int_E \nu(A_k - x) d\mu(x) = \int_E \sum_{k=1}^n c_k \nu(A_k - x) d\mu(x) \\ &= \int_E \sum_{k=1}^n c_k \mathcal{X}_{A_k - x}(y) d\nu(y) d\mu(x) = \int_E \int_E s(x+y) d\nu(y) d\mu(x) \end{aligned}$$

Approximating with simple functions it follows that whenever f is bounded and measurable or nonnegative and measurable,

$$\int_E f d(\mu * \nu) = \int_E \int_E f(x+y) d\nu(y) d\mu(x) \quad (27.14)$$

Therefore, letting $Z = X + Y$, and λ the distribution of Z , it follows from independence of X and Y that for $t^* \in E'$,

$$\phi_\lambda(t^*) \equiv E\left(e^{it^*(Z)}\right) = E\left(e^{it^*(X+Y)}\right) = E\left(e^{it^*(X)}\right) E\left(e^{it^*(Y)}\right)$$

But also, it follows from 27.14

$$\begin{aligned} \phi_{(\mu * \nu)}(t^*) &= \int_E e^{it^*(z)} d(\mu * \nu)(z) = \int_E \int_E e^{it^*(x+y)} d\nu(y) d\mu(x) \\ &= \int_E \int_E e^{it^*(x)} e^{it^*(y)} d\nu(y) d\mu(x) \\ &= \left(\int_E e^{it^*(y)} d\nu(y) \right) \left(\int_E e^{it^*(x)} d\mu(x) \right) \\ &= E\left(e^{it^*(X)}\right) E\left(e^{it^*(Y)}\right) \end{aligned}$$

Since $\phi_\lambda(t^*) = \phi_{(\mu * \nu)}(t^*)$, it follows $\lambda = \mu * \nu$.

Note the last part of this argument shows the characteristic function of a convolution equals the product of the characteristic functions. ■

Chapter 28

The Normal Distribution

This particular distribution is likely the most important one in statistics and it will be essential to understand in developing the Wiener process later. To begin with, $\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-\frac{1}{2}u^2} du = 1$ as is easily shown as done earlier by the standard calculus trick of

$$I = \int_{\mathbb{R}} e^{-\frac{1}{2}u^2} du, I^2 = \int_{\mathbb{R}} \int_{\mathbb{R}} e^{-\frac{1}{2}(u^2+v^2)} dudv$$

and then changing to polar coordinates to obtain $I^2 = 2\pi$. I will use this identity whenever convenient. Also useful is the following lemma.

Lemma 28.0.1 $\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-\frac{1}{2}(u-it)^2} du = 1$.

Proof: $e^{-\frac{1}{2}(u-it)^2} = e^{-\frac{1}{2}(u^2-2itu-t^2)} = e^{-\frac{1}{2}u^2} e^{\frac{1}{2}t^2} (\cos(tu) + i \sin(tu))$ and so, the integral equals

$$\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-\frac{1}{2}u^2} e^{\frac{1}{2}t^2} \cos(tu) du$$

Now let $f(t) \equiv \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-\frac{1}{2}u^2} e^{\frac{1}{2}t^2} \cos(tu) du$. Using the dominated convergence theorem,

$$\begin{aligned} f'(t) &= \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \frac{d}{dt} \left(e^{-\frac{1}{2}u^2} e^{\frac{1}{2}t^2} \cos(tu) \right) du \\ &= \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \left(e^{-\frac{1}{2}u^2} \left(t e^{\frac{1}{2}t^2} \cos(tu) - e^{\frac{1}{2}t^2} u \sin(tu) \right) \right) du \end{aligned}$$

Now $f(0)$ is known to be 1. Assume then that $t \neq 0$.

$$-\frac{1}{\sqrt{2\pi}} e^{\frac{1}{2}t^2} \int_{\mathbb{R}} e^{-\frac{1}{2}u^2} u \sin(tu) du = \frac{1}{\sqrt{2\pi}} e^{\frac{1}{2}t^2} \int_{\mathbb{R}} e^{-\frac{1}{2}u^2} t \cos(tu) du$$

and this shows that $f'(t) = 0$ so $f(t)$ is the constant 1. ■

28.1 The Multivariate Normal Distribution

The multivariate normal distribution is very important in statistics and it will be shown in this chapter why this is the case.

Definition 28.1.1 A random vector \mathbf{X} , with values in \mathbb{R}^p has a multivariate normal distribution written as $\mathbf{X} \sim N_p(\mathbf{m}, \Sigma)$ if for all Borel $E \subseteq \mathbb{R}^p$,

$$\lambda_{\mathbf{X}}(E) = \int_{\mathbb{R}^p} \mathcal{X}_E(\mathbf{x}) \frac{1}{(2\pi)^{p/2} \det(\Sigma)^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\mathbf{m})^* \Sigma^{-1}(\mathbf{x}-\mathbf{m})} d\mathbf{x}$$

for μ a given vector and Σ a given positive definite symmetric matrix, called the covariance matrix. In case $p = 1$, this is called the variance.

Theorem 28.1.2 For $\mathbf{X} \sim N_p(\mathbf{m}, \Sigma)$, $\mathbf{m} = E(\mathbf{X})$ and

$$\Sigma = E((\mathbf{X} - \mathbf{m})(\mathbf{X} - \mathbf{m})^*).$$

Proof: Let R be an orthogonal transformation such that

$$R\Sigma R^* = D = \text{diag}(\sigma_1^2, \dots, \sigma_p^2), \quad \sigma_i > 0.$$

Changing the variable by $\mathbf{x} - \mathbf{m} = R^* \mathbf{y}$,

$$\begin{aligned} E(\mathbf{X}) &\equiv \int_{\mathbb{R}^p} \mathbf{x} e^{-\frac{1}{2}(\mathbf{x}-\mathbf{m})^* \Sigma^{-1}(\mathbf{x}-\mathbf{m})} d\mathbf{x} \left(\frac{1}{(2\pi)^{p/2} \det(\Sigma)^{1/2}} \right) \\ &= \int_{\mathbb{R}^p} (R^* \mathbf{y} + \mathbf{m}) e^{-\frac{1}{2} \mathbf{y}^* D^{-1} \mathbf{y}} d\mathbf{y} \left(\frac{1}{(2\pi)^{p/2} \prod_{i=1}^p \sigma_i} \right) \\ &= \mathbf{m} \int_{\mathbb{R}^p} e^{-\frac{1}{2} \mathbf{y}^* D^{-1} \mathbf{y}} d\mathbf{y} \left(\frac{1}{(2\pi)^{p/2} \prod_{i=1}^p \sigma_i} \right) = \mathbf{m} \end{aligned}$$

by Fubini's theorem and the easy to establish formula $\frac{1}{\sqrt{2\pi}\sigma} \int_{\mathbb{R}} e^{-\frac{y^2}{2\sigma^2}} dy = 1$, (let $u = y/\sigma$),

Next let

$$M \equiv E((\mathbf{X} - \mathbf{m})(\mathbf{X} - \mathbf{m})^*)$$

Thus, changing the variable as above by $\mathbf{x} - \mathbf{m} = R^* \mathbf{y}$

$$\begin{aligned} M &= \int_{\mathbb{R}^p} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^* e^{-\frac{1}{2}(\mathbf{x}-\mathbf{m})^* \Sigma^{-1}(\mathbf{x}-\mathbf{m})} d\mathbf{x} \left(\frac{1}{(2\pi)^{p/2} \det(\Sigma)^{1/2}} \right) \\ &= R^* \int_{\mathbb{R}^p} \mathbf{y} \mathbf{y}^* e^{-\frac{1}{2} \mathbf{y}^* D^{-1} \mathbf{y}} d\mathbf{y} \left(\frac{1}{(2\pi)^{p/2} \prod_{j=1}^p \sigma_j} \right) R \end{aligned}$$

If $i \neq j$, $(RMR^*)_{ij} = \int_{\mathbb{R}^p} y_i y_j e^{-\frac{1}{2} \mathbf{y}^* D^{-1} \mathbf{y}} d\mathbf{y} \left(\frac{1}{(2\pi)^{p/2} \prod_{k=1}^p \sigma_k} \right) = 0$ so RMR^* is a diagonal matrix.

$$(RMR^*)_{ii} = \int_{\mathbb{R}^p} y_i^2 e^{-\frac{1}{2} \mathbf{y}^* D^{-1} \mathbf{y}} d\mathbf{y} \left(\frac{1}{(2\pi)^{p/2} \prod_{j=1}^p \sigma_j} \right).$$

Using Fubini's theorem and the easy to establish equations,

$$\frac{1}{\sqrt{2\pi}\sigma} \int_{\mathbb{R}} e^{-\frac{y^2}{2\sigma^2}} dy = 1, \quad \frac{1}{\sqrt{2\pi}\sigma} \int_{\mathbb{R}} y^2 e^{-\frac{y^2}{2\sigma^2}} dy = \sigma^2,$$

it follows $(RMR^*)_{ii} = \sigma_i^2$. Hence $RMR^* = D$ and so $M = R^*DR = \Sigma$. ■

Theorem 28.1.3 Suppose $\mathbf{X}_1 \sim N_p(\mathbf{m}_1, \Sigma_1)$, $\mathbf{X}_2 \sim N_p(\mathbf{m}_2, \Sigma_2)$ and the two random vectors are independent. Then

$$\mathbf{X}_1 + \mathbf{X}_2 \sim N_p(\mathbf{m}_1 + \mathbf{m}_2, \Sigma_1 + \Sigma_2). \quad (28.1)$$

Also, if $\mathbf{X} \sim N_p(\mathbf{m}, \Sigma)$ then $-\mathbf{X} \sim N_p(-\mathbf{m}, \Sigma)$. Furthermore, if $\mathbf{X} \sim N_p(\mathbf{m}, \Sigma)$ then

$$E(e^{it \cdot \mathbf{X}}) = e^{it \cdot \mathbf{m}} e^{-\frac{1}{2} t^* \Sigma t} \quad (28.2)$$

If a is a constant and $\mathbf{X} \sim N_p(\mathbf{m}, \Sigma)$, then $a\mathbf{X} \sim N_p(a\mathbf{m}, a^2\Sigma)$.

Proof: Consider $E(e^{it \cdot X})$ for $X \sim N_p(\mathbf{m}, \Sigma)$.

$$E(e^{it \cdot X}) \equiv \frac{1}{(2\pi)^{p/2} (\det \Sigma)^{1/2}} \int_{\mathbb{R}^p} e^{it \cdot x} e^{-\frac{1}{2}(x-\mathbf{m})^* \Sigma^{-1} (x-\mathbf{m})} dx.$$

Let R be an orthogonal transformation such that

$$R \Sigma R^* = D = \text{diag}(\sigma_1^2, \dots, \sigma_p^2).$$

Let $R(x - \mathbf{m}) = \mathbf{y}$. Then

$$E(e^{it \cdot X}) = \frac{1}{(2\pi)^{p/2} \prod_{i=1}^p \sigma_i} \int_{\mathbb{R}^p} e^{it \cdot (R^* \mathbf{y} + \mathbf{m})} e^{-\frac{1}{2} \mathbf{y}^* D^{-1} \mathbf{y}} d\mathbf{y}.$$

Therefore

$$E(e^{it \cdot X}) = \frac{1}{(2\pi)^{p/2} \prod_{i=1}^p \sigma_i} \int_{\mathbb{R}^p} e^{is \cdot (\mathbf{y} + R\mathbf{m})} e^{-\frac{1}{2} \mathbf{y}^* D^{-1} \mathbf{y}} d\mathbf{y}$$

where $s = Rt$. This equals

$$\begin{aligned} & e^{it \cdot \mathbf{m}} \prod_{i=1}^p \left(\int_{\mathbb{R}} e^{is_i y_i} e^{-\frac{1}{2\sigma_i^2} y_i^2} dy_i \right) \frac{1}{\sqrt{2\pi} \sigma_i} \\ &= e^{it \cdot \mathbf{m}} \prod_{i=1}^p \left(\int_{\mathbb{R}} e^{is_i \sigma_i u} e^{-\frac{1}{2} u^2} du \right) \frac{1}{\sqrt{2\pi}} \\ &= e^{it \cdot \mathbf{m}} \prod_{i=1}^p e^{-\frac{1}{2} s_i^2 \sigma_i^2} \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-\frac{1}{2} (u - is_i \sigma_i)^2} du \end{aligned}$$

By Lemma 28.0.1, this equals $e^{it \cdot \mathbf{m}} e^{-\frac{1}{2} \Sigma_{i=1}^p s_i^2 \sigma_i^2} = e^{it \cdot \mathbf{m}} e^{-\frac{1}{2} t^* \Sigma t}$. This proves 28.2.

Since X_1 and X_2 are independent, $e^{it \cdot X_1}$ and $e^{it \cdot X_2}$ are also independent. Hence

$$E(e^{it \cdot X_1 + X_2}) = E(e^{it \cdot X_1}) E(e^{it \cdot X_2}).$$

Thus,

$$\begin{aligned} E(e^{it \cdot X_1 + X_2}) &= E(e^{it \cdot X_1}) E(e^{it \cdot X_2}) = e^{it \cdot \mathbf{m}_1} e^{-\frac{1}{2} t^* \Sigma_1 t} e^{it \cdot \mathbf{m}_2} e^{-\frac{1}{2} t^* \Sigma_2 t} \\ &= e^{it \cdot (\mathbf{m}_1 + \mathbf{m}_2)} e^{-\frac{1}{2} t^* (\Sigma_1 + \Sigma_2) t}, \end{aligned}$$

which, as shown above is the characteristic function of a random vector distributed as $N_p(\mathbf{m}_1 + \mathbf{m}_2, \Sigma_1 + \Sigma_2)$. Now it follows that $X_1 + X_2 \sim N_p(\mathbf{m}_1 + \mathbf{m}_2, \Sigma_1 + \Sigma_2)$ by Theorem 27.1.4. This proves 28.1.

The assertion about $-X$ is also easy to see because

$$\begin{aligned} E(e^{it \cdot (-X)}) &= E(e^{i(-t) \cdot X}) \\ &= \frac{1}{(2\pi)^{p/2} (\det \Sigma)^{1/2}} \int_{\mathbb{R}^p} e^{i(-t) \cdot x} e^{-\frac{1}{2}(x-\mathbf{m})^* \Sigma^{-1} (x-\mathbf{m})} dx \\ &= \frac{1}{(2\pi)^{p/2} (\det \Sigma)^{1/2}} \int_{\mathbb{R}^p} e^{it \cdot x} e^{-\frac{1}{2}(x+\mathbf{m})^* \Sigma^{-1} (x+\mathbf{m})} dx \end{aligned}$$

which is the characteristic function of a random variable which is $N(-\mathbf{m}, \Sigma)$. Theorem 27.1.4 again implies $-\mathbf{X} \sim N(-\mathbf{m}, \Sigma)$. Finally consider the last claim. You apply what is known about \mathbf{X} with \mathbf{t} replaced with $a\mathbf{t}$ and then massage things. This gives the characteristic function for $a\mathbf{X}$ is given by

$$E(\exp(i\mathbf{t} \cdot a\mathbf{X})) = \exp(i\mathbf{t} \cdot a\mathbf{m}) \exp\left(-\frac{1}{2}\mathbf{t}^* \Sigma a^2 \mathbf{t}\right)$$

which is the characteristic function of a normal random vector having mean $a\mathbf{m}$ and covariance $a^2\Sigma$. ■

28.2 Linear Combinations

Following [44] a random vector has a generalized normal distribution if its characteristic function is given as $e^{it \cdot \mathbf{m}} e^{-\frac{1}{2} \mathbf{t}^* \Sigma \mathbf{t}}$ where Σ is symmetric and has nonnegative eigenvalues. For a random real valued variable, \mathbf{m} is scalar and so is Σ so the characteristic function of such a generalized normally distributed random variable is $e^{it\mu} e^{-\frac{1}{2} t^2 \sigma^2}$. These generalized normal distributions do not require Σ to be invertible, only that the eigenvalues be nonnegative. In one dimension this would correspond the characteristic function of a dirac measure having point mass 1 at μ . In higher dimensions, it could be a mixture of such things with more familiar things. I won't try very hard to distinguish between generalized normal distributions and normal distributions in which the covariance matrix has all positive eigenvalues. These generalized normal distributions are discussed more a little later.

Here are some other interesting results about normal distributions found in [44]. The next theorem has to do with the question whether a random vector is normally distributed in the above generalized sense.

Theorem 28.2.1 *Let $\mathbf{X} = (X_1, \dots, X_p)$ where each X_i is a real valued random variable. Then \mathbf{X} is normally distributed in the above generalized sense if and only if every linear combination, $\sum_{j=1}^p a_j X_j$ is normally distributed. In this case the mean of \mathbf{X} is*

$$\mathbf{m} = (E(X_1), \dots, E(X_p))$$

and the covariance matrix for \mathbf{X} is

$$\Sigma_{jk} = E((X_j - m_j)(X_k - m_k)^*).$$

Proof: Suppose first \mathbf{X} is normally distributed. Then its characteristic function is of the form

$$\phi_{\mathbf{X}}(\mathbf{t}) = E(e^{it \cdot \mathbf{X}}) = e^{it \cdot \mathbf{m}} e^{-\frac{1}{2} \mathbf{t}^* \Sigma \mathbf{t}}.$$

Then letting $\mathbf{a} = (a_1, \dots, a_p)$

$$E\left(e^{it \sum_{j=1}^p a_j X_j}\right) = E(e^{it \mathbf{a} \cdot \mathbf{X}}) = e^{it \mathbf{a} \cdot \mathbf{m}} e^{-\frac{1}{2} \mathbf{a}^* \Sigma \mathbf{a} t^2}$$

which is the characteristic function of a normally distributed random variable with mean $\mathbf{a} \cdot \mathbf{m}$ and variance $\sigma^2 = \mathbf{a}^* \Sigma \mathbf{a}$. This proves half of the theorem. If \mathbf{X} is normally distributed, then every linear combination is normally distributed.

Conversely, suppose every linear combination is normally distributed. Next suppose $\sum_{j=1}^p a_j X_j = \mathbf{a} \cdot \mathbf{X}$ is normally distributed with mean μ and variance σ^2 so that its characteristic function is given as $e^{i\mu} e^{-\frac{1}{2}t^2\sigma^2}$. I will now relate μ and σ^2 to various quantities involving the X_j . Letting $m_j = E(X_j)$, $\mathbf{m} = (m_1, \dots, m_p)^*$

$$\begin{aligned}\mu &= \sum_{j=1}^p a_j E(X_j) = \sum_{j=1}^p a_j m_j, \quad \sigma^2 = E \left(\left(\sum_{j=1}^p a_j X_j - \sum_{j=1}^p a_j m_j \right)^2 \right) \\ &= E \left(\left(\sum_{j=1}^p a_j (X_j - m_j) \right)^2 \right) = \sum_{j,k} a_j a_k E((X_j - m_j)(X_k - m_k))\end{aligned}$$

It follows the mean of the random variable, $\mathbf{a} \cdot \mathbf{X}$ is $\mu = \sum_j a_j m_j = \mathbf{a} \cdot \mathbf{m}$ and its variance is

$$\sigma^2 = \mathbf{a}^* E((\mathbf{X} - \mathbf{m})(\mathbf{X} - \mathbf{m})^*) \mathbf{a}$$

Therefore, $E(e^{i\mathbf{a} \cdot \mathbf{X}}) = e^{i\mu} e^{-\frac{1}{2}t^2\sigma^2} = e^{i\mathbf{a} \cdot \mathbf{m}} e^{-\frac{1}{2}t^2 \mathbf{a}^* E((\mathbf{X} - \mathbf{m})(\mathbf{X} - \mathbf{m})^*) \mathbf{a}}$. Letting $\mathbf{s} = t\mathbf{a}$ this shows

$$E(e^{i\mathbf{s} \cdot \mathbf{X}}) = e^{i\mathbf{s} \cdot \mathbf{m}} e^{-\frac{1}{2}\mathbf{s}^* E((\mathbf{X} - \mathbf{m})(\mathbf{X} - \mathbf{m})^*) \mathbf{s}} = e^{i\mathbf{s} \cdot \mathbf{m}} e^{-\frac{1}{2}\mathbf{s}^* \Sigma \mathbf{s}}$$

which is the characteristic function of a normally distributed random variable with \mathbf{m} given above and Σ given by

$$\Sigma_{jk} = E((X_j - m_j)(X_k - m_k)).$$

By assumption, \mathbf{a} is completely arbitrary and so it follows that \mathbf{s} is also. Hence, \mathbf{X} is normally distributed as claimed. ■

Corollary 28.2.2 Let $\mathbf{X} = (X_1, \dots, X_p)$, $\mathbf{Y} = (Y_1, \dots, Y_p)$ where each X_i, Y_i is a real valued random variable. Suppose also that for every $\mathbf{a} \in \mathbb{R}^p$, $\mathbf{a} \cdot \mathbf{X}$ and $\mathbf{a} \cdot \mathbf{Y}$ are both normally distributed with the same mean and variance. Then \mathbf{X} and \mathbf{Y} are both multivariate normal random vectors with the same mean and variance.

Proof: In the Proof of Theorem 28.2.1 the proof implies that the characteristic functions of $\mathbf{a} \cdot \mathbf{X}$ and $\mathbf{a} \cdot \mathbf{Y}$ are both of the form $e^{i\mathbf{a} \cdot \mathbf{m}} e^{-\frac{1}{2}\sigma^2 t^2}$. Then as in the proof of that theorem, it must be the case that $m = \sum_{j=1}^p a_j m_j$ where $E(X_i) = m_i = E(Y_i)$ and

$$\sigma^2 = \mathbf{a}^* E((\mathbf{X} - \mathbf{m})(\mathbf{X} - \mathbf{m})^*) \mathbf{a} = \mathbf{a}^* E((\mathbf{Y} - \mathbf{m})(\mathbf{Y} - \mathbf{m})^*) \mathbf{a}$$

and this last equation must hold for every \mathbf{a} . Therefore,

$$E((\mathbf{X} - \mathbf{m})(\mathbf{X} - \mathbf{m})^*) = E((\mathbf{Y} - \mathbf{m})(\mathbf{Y} - \mathbf{m})^*) \equiv \Sigma$$

and so the characteristic function of both \mathbf{X} and \mathbf{Y} is $e^{i\mathbf{s} \cdot \mathbf{m}} e^{-\frac{1}{2}\mathbf{s}^* \Sigma \mathbf{s}}$ as in the proof of Theorem 28.2.1. ■

Theorem 28.2.3 Suppose $\mathbf{X} = (X_1, \dots, X_p)$ is normally distributed with mean \mathbf{m} and covariance Σ . Then if X_1 is uncorrelated with any of the X_i , meaning

$$E((X_1 - m_1)(X_j - m_j)) = 0 \text{ for } j > 1,$$

then X_1 and (X_2, \dots, X_p) are both normally distributed and the two random vectors are independent. Here $m_j \equiv E(X_j)$. More generally, if the covariance matrix is a diagonal matrix, the random variables, $\{X_1, \dots, X_p\}$ are independent.

Proof: From Theorem 28.1.2 $\Sigma = E((\mathbf{X} - \mathbf{m})(\mathbf{X} - \mathbf{m})^*)$. Then by assumption,

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \mathbf{0} \\ \mathbf{0} & \Sigma_{p-1} \end{pmatrix}. \quad (28.3)$$

I need to verify that if $E \in \sigma(X_1)$ and $F \in \sigma(X_2, \dots, X_p)$, then $P(E \cap F) = P(E)P(F)$.

Let $E = X_1^{-1}(A)$ and $F = (X_2, \dots, X_p)^{-1}(B)$ where A and B are Borel sets in \mathbb{R} and \mathbb{R}^{p-1} respectively. Thus I need to verify that

$$\begin{aligned} P([(X_1, (X_2, \dots, X_p)) \in (A, B)]) &= \\ \mu_{(X_1, (X_2, \dots, X_p))}(A \times B) &= \mu_{X_1}(A) \mu_{(X_2, \dots, X_p)}(B). \end{aligned} \quad (28.4)$$

Using 28.3, Fubini's theorem, and definitions,

$$\begin{aligned} \mu_{(X_1, (X_2, \dots, X_p))}(A \times B) &= \\ \int_{\mathbb{R}^p} \mathcal{X}_{A \times B}(\mathbf{x}) \frac{1}{(2\pi)^{p/2} \det(\Sigma)^{1/2}} e^{-\frac{1}{2}(\mathbf{x} - \mathbf{m})^* \Sigma^{-1}(\mathbf{x} - \mathbf{m})} d\mathbf{x} &= \\ = \int_{\mathbb{R}} \mathcal{X}_A(x_1) \int_{\mathbb{R}^{p-1}} \mathcal{X}_B(X_2, \dots, X_p) \cdot & \\ \frac{1}{(2\pi)^{(p-1)/2} \sqrt{2\pi} (\sigma_1^2)^{1/2} \det(\Sigma_{p-1})^{1/2}} e^{-\frac{(x_1 - m_1)^2}{2\sigma_1^2}} \cdot & \\ e^{-\frac{1}{2}(\mathbf{x}' - \mathbf{m}')^* \Sigma_{p-1}^{-1}(\mathbf{x}' - \mathbf{m}')} d\mathbf{x}' dx_1 & \end{aligned}$$

where $\mathbf{x}' = (x_2, \dots, x_p)$ and $\mathbf{m}' = (m_2, \dots, m_p)$. Now this equals

$$\begin{aligned} \int_{\mathbb{R}} \mathcal{X}_A(x_1) \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{(x_1 - m_1)^2}{2\sigma_1^2}} \cdot & \\ \int_B \frac{1}{(2\pi)^{(p-1)/2} \det(\Sigma_{p-1})^{1/2}} e^{-\frac{1}{2}(\mathbf{x}' - \mathbf{m}')^* \Sigma_{p-1}^{-1}(\mathbf{x}' - \mathbf{m}')} d\mathbf{x}' dx_1. & \end{aligned} \quad (28.5)$$

In case $B = \mathbb{R}^{p-1}$, the inside integral equals 1 and

$$\mu_{X_1}(A) = \mu_{(X_1, (X_2, \dots, X_p))}(A \times \mathbb{R}^{p-1}) = \int_{\mathbb{R}} \mathcal{X}_A(x_1) \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{(x_1 - m_1)^2}{2\sigma_1^2}} dx_1$$

which shows X_1 is normally distributed as claimed. Similarly, letting $A = \mathbb{R}$,

$$\begin{aligned} \mu_{(X_2, \dots, X_p)}(B) &= \mu_{(X_1, (X_2, \dots, X_p))}(\mathbb{R} \times B) \\ = \int_B \frac{1}{(2\pi)^{(p-1)/2} \det(\Sigma_{p-1})^{1/2}} e^{-\frac{1}{2}(\mathbf{x}' - \mathbf{m}')^* \Sigma_{p-1}^{-1}(\mathbf{x}' - \mathbf{m}')} d\mathbf{x}' & \end{aligned}$$

and (X_2, \dots, X_p) is also normally distributed with mean \mathbf{m}' and covariance Σ_{p-1} . Now from 28.5, 28.4 follows. In case the covariance matrix is diagonal, the above reasoning extends in an obvious way to prove the random variables, $\{X_1, \dots, X_p\}$ are independent.

However, another way to prove this is to use Proposition 27.4.1 on Page 744 and consider the characteristic function. Let $E(X_j) = m_j$ and $P = \sum_{j=1}^p t_j X_j$. Then since \mathbf{X} is normally distributed and the covariance is a diagonal,

$$D \equiv \begin{pmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_p^2 \end{pmatrix}$$

$$\begin{aligned} E(e^{iP}) &= E(e^{it \cdot \mathbf{X}}) = e^{it \cdot \mathbf{m}} e^{-\frac{1}{2} \mathbf{t}^* \Sigma \mathbf{t}} = \exp \left(\sum_{j=1}^p it_j m_j - \frac{1}{2} t_j^2 \sigma_j^2 \right) \\ &= \prod_{j=1}^p \exp \left(it_j m_j - \frac{1}{2} t_j^2 \sigma_j^2 \right) \end{aligned} \quad (28.6)$$

Also,

$$E(e^{it_j X_j}) = E \left(\exp \left(it_j X_j + \sum_{k \neq j} i0 X_k \right) \right) = \exp \left(it_j m_j - \frac{1}{2} t_j^2 \sigma_j^2 \right)$$

With 28.6, this shows $E(e^{iP}) = \prod_{j=1}^p E(e^{it_j X_j})$ which shows by Proposition 27.4.1 that the random variables, $\{X_1, \dots, X_p\}$ are independent. ■

28.3 Finding Moments

Let X be a random variable with characteristic function

$$\phi_X(t) \equiv E(\exp(itX))$$

Then this can be used to find moments of the random variable assuming they exist. The k^{th} moment is defined as $E(X^k)$. This can be done by using the dominated convergence theorem to differentiate the characteristic function with respect to t and then plugging in $t = 0$. For example, $\phi'_X(t) = E(iX \exp(itX))$ and now plugging in $t = 0$ you get $iE(X)$. Doing another differentiation you obtain $\phi''_X(t) = E(-X^2 \exp(itX))$ and plugging in $t = 0$ you get $-E(X^2)$ and so forth.

An important case is where X is normally distributed with mean 0 and variance σ^2 . In this case, as shown above, the characteristic function is $e^{-\frac{1}{2}t^2\sigma^2}$. Also all moments exist when X is normally distributed. So what are these moments? $D_t(e^{-\frac{1}{2}t^2\sigma^2}) = -t\sigma^2 e^{-\frac{1}{2}t^2\sigma^2}$ and plugging in $t = 0$ you find the mean equals 0 as expected.

$$D_t(-t\sigma^2 e^{-\frac{1}{2}t^2\sigma^2}) = -\sigma^2 e^{-\frac{1}{2}t^2\sigma^2} + t^2\sigma^4 e^{-\frac{1}{2}t^2\sigma^2}$$

and plugging in $t = 0$ you find the second moment is σ^2 . Then do it again.

$$D_t(-\sigma^2 e^{-\frac{1}{2}t^2\sigma^2} + t^2\sigma^4 e^{-\frac{1}{2}t^2\sigma^2}) = 3\sigma^4 t e^{-\frac{1}{2}t^2\sigma^2} - t^3\sigma^6 e^{-\frac{1}{2}t^2\sigma^2}$$

Then $E(X^3) = 0$.

$$D_t(3\sigma^4 t e^{-\frac{1}{2}t^2\sigma^2} - t^3\sigma^6 e^{-\frac{1}{2}t^2\sigma^2}) = 3\sigma^4 e^{-\frac{1}{2}t^2\sigma^2} - 6\sigma^6 t^2 e^{-\frac{1}{2}t^2\sigma^2} + t^4\sigma^8 e^{-\frac{1}{2}t^2\sigma^2}$$

and so $E(X^4) = 3\sigma^4$. By now you can see the pattern. If you continue this way, you find the odd moments are all 0 and

$$E(X^{2m}) = C_m (\sigma^2)^m. \quad (28.7)$$

This is an important observation. In the case of \mathbf{X} a random vector, you have $\phi_{\mathbf{X}}(\mathbf{t}) \equiv E(\exp(it \cdot \mathbf{X}))$ and by taking $\frac{d}{dt_j}$, you can follow the above procedure to obtain $E(X_j)$ and then by using successive differentiations, you can find $E(X_i^n)$ or any polynomial in the X_i assuming the expectations exist.

28.4 Prokhorov and Levy Theorems

Recall one can define the characteristic function of a probability measure μ as $\int_{\mathbb{R}^p} e^{it \cdot x} d\mu$. In a sense it is more natural. One can also generalize to replace \mathbb{R}^p with E a Banach space in which the dot product $t \cdot x$ is replaced with $t(x)$ where $t \in E'$. However, the main interest here is in \mathbb{R}^p .

Definition 28.4.1 A set of Borel probability measures $\{\mu_n\}_{n=1}^\infty$ defined on a Polish space E is called “tight” if for all $\varepsilon > 0$ there exists a compact set, K_ε such that

$$\mu_n([x \notin K_\varepsilon]) < \varepsilon$$

for all μ_n .

How do you determine in general that a set of probability measures is tight?

Lemma 28.4.2 Let E be a separable complete metric space and let Λ be a set of Borel probability measures. Then Λ is tight if and only if for every $\varepsilon > 0$ and $r > 0$ there exists a finite collection of balls, $\{B(a_i, r)\}_{i=1}^m$ such that

$$\mu\left(\bigcup_{i=1}^m \overline{B(a_i, r)}\right) > 1 - \varepsilon$$

for every $\mu \in \Lambda$.

Proof: If Λ is tight, then there exists a compact set, K_ε such that

$$\mu(K_\varepsilon) > 1 - \varepsilon$$

for all $\mu \in \Lambda$. Then consider the open cover, $\{B(x, r) : x \in K_\varepsilon\}$. Finitely many of these cover K_ε and this yields the above condition.

Now suppose the above condition and let

$$C_n \equiv \bigcup_{i=1}^{m_n} \overline{B(a_i^n, 1/n)}$$

satisfy $\mu(C_n) > 1 - \varepsilon/2^n$ for all $\mu \in \Lambda$. Then let $K_\varepsilon \equiv \bigcap_{n=1}^\infty C_n$. This set K_ε is a compact set because it is a closed subset of a complete metric space and is therefore complete, and it is also totally bounded by construction. For $\mu \in \Lambda$,

$$\mu(K_\varepsilon^C) = \mu\left(\bigcup_n C_n^C\right) \leq \sum_n \mu(C_n^C) < \sum_n \frac{\varepsilon}{2^n} = \varepsilon$$

Therefore, Λ is tight. ■

In case the Polish space is \mathbb{R}^p , the case of most interest, there is a very nice condition in terms of characteristic functions which gives “tightness”.

Lemma 28.4.3 *If $\{\mu_n\}$ is a sequence of Borel probability measures defined on the Borel sets of \mathbb{R}^p such that*

$$\lim_{n \rightarrow \infty} \phi_{\mu_n}(\mathbf{t}) = \psi(\mathbf{t})$$

for all \mathbf{t} , where $\psi(\mathbf{0}) = 1$ and ψ is continuous at $\mathbf{0}$, then $\{\mu_n\}_{n=1}^\infty$ is tight.

Proof: Let \mathbf{e}_j be the j^{th} standard unit basis vector. Letting $\mathbf{t} = t\mathbf{e}_j$ in the definition and $u > 0$

$$\begin{aligned} \left| \frac{1}{u} \int_{-u}^u (1 - \phi_{\mu_n}(t\mathbf{e}_j)) dt \right| &= \left| \frac{1}{u} \int_{-u}^u \left(1 - \int_{\mathbb{R}^p} e^{itx_j} d\mu_n(x) \right) dt \right| \\ &= \left| \frac{1}{u} \int_{-u}^u \left(\int_{\mathbb{R}^p} (1 - e^{itx_j}) d\mu_n(x) \right) dt \right| = \left| \int_{\mathbb{R}^p} \frac{1}{u} \int_{-u}^u (1 - e^{itx_j}) dt d\mu_n(x) \right| \\ &= \left| 2 \int_{\mathbb{R}^p} \left(1 - \frac{\sin(ux_j)}{ux_j} \right) d\mu_n(x) \right| \geq 2 \int_{[|x_j| \geq \frac{2}{u}]} \left(1 - \frac{1}{|ux_j|} \right) d\mu_n(x) \\ &\geq 2 \int_{[|x_j| \geq \frac{2}{u}]} \left(1 - \frac{1}{u(2/u)} \right) d\mu_n(x) = \int_{[|x_j| \geq \frac{2}{u}]} 1 d\mu_n(x) = \mu_n \left(\left[\mathbf{x} : |x_j| \geq \frac{2}{u} \right] \right). \end{aligned} \quad (28.8)$$

If $\varepsilon > 0$ is given, there exists $r > 0$ such that if $u \leq r$, $\frac{1}{u} \int_{-u}^u (1 - \psi(t\mathbf{e}_j)) dt < \varepsilon/p$ for all $j = 1, \dots, p$ and so, by the dominated convergence theorem, the same is true with ϕ_{μ_n} in place of ψ provided n is large enough, say $n \geq N(r)$. Thus, from 28.8, if $n \geq N(r)$, $\mu_n([x : |x_j| > 2r]) < \varepsilon/p$ for all $j \in \{1, \dots, p\}$. It follows for $n \geq N(r)$,

$$\mu_n([x : \|x\|_\infty > 2r]) < \varepsilon.$$

and so let $K_\varepsilon \equiv [-r, r]^p$. ■

In the case of \mathbb{R}^p , and μ_n, μ Borel probability measures, convergence of characteristic functions yields something interesting for $\psi \in \mathcal{G}$ or \mathfrak{S} , the Schwartz class.

Lemma 28.4.4 *If $\phi_{\mu_n}(\mathbf{t}) \rightarrow \phi_\mu(\mathbf{t})$ for all \mathbf{t} , then whenever $\psi \in \mathfrak{S}$, the Schwartz class,*

$$\mu_n(\psi) \equiv \int_{\mathbb{R}^p} \psi(\mathbf{y}) d\mu_n(\mathbf{y}) \rightarrow \int_{\mathbb{R}^p} \psi(\mathbf{y}) d\mu(\mathbf{y}) \equiv \mu(\psi)$$

as $n \rightarrow \infty$.

Proof: By definition, $\phi_\mu(\mathbf{y}) \equiv \int_{\mathbb{R}^p} e^{i\mathbf{y} \cdot \mathbf{x}} d\mu(\mathbf{x})$. Also remember the inverse Fourier transform. Letting $\psi \in \mathfrak{S}$, the Schwartz class,

$$\begin{aligned} F^{-1}(\mu)(\psi) &\equiv \mu(F^{-1}\psi) \equiv \int_{\mathbb{R}^p} F^{-1}\psi d\mu \\ &= \frac{1}{(2\pi)^{p/2}} \int_{\mathbb{R}^p} \int_{\mathbb{R}^p} e^{i\mathbf{y} \cdot \mathbf{x}} \psi(\mathbf{x}) d\mathbf{x} d\mu(\mathbf{y}) \\ &= \frac{1}{(2\pi)^{p/2}} \int_{\mathbb{R}^p} \psi(\mathbf{x}) \int_{\mathbb{R}^p} e^{i\mathbf{y} \cdot \mathbf{x}} d\mu(\mathbf{y}) d\mathbf{x} = \frac{1}{(2\pi)^{p/2}} \int_{\mathbb{R}^p} \psi(\mathbf{x}) \phi_\mu(\mathbf{x}) d\mathbf{x} \end{aligned}$$

and so, considered as elements of \mathfrak{S}^* or \mathcal{G}^* , $F^{-1}(\mu) = \phi_\mu(\cdot) (2\pi)^{-(p/2)} \in L^\infty$. By the dominated convergence theorem

$$\begin{aligned} (2\pi)^{p/2} F^{-1}(\mu_n)(\psi) &\equiv \int_{\mathbb{R}^p} \phi_{\mu_n}(\mathbf{t}) \psi(\mathbf{t}) d\mathbf{t} \rightarrow \int_{\mathbb{R}^p} \phi_\mu(\mathbf{t}) \psi(\mathbf{t}) d\mathbf{t} \\ &= (2\pi)^{p/2} F^{-1}(\mu)(\psi) \end{aligned}$$

whenever $\psi \in \mathfrak{S}$ or \mathcal{G} . Thus

$$\begin{aligned}\mu_n(\psi) &= FF^{-1}\mu_n(\psi) \equiv F^{-1}\mu_n(F\psi) \rightarrow F^{-1}\mu(F\psi) \\ &\equiv F^{-1}F\mu(\psi) = \mu(\psi). \blacksquare\end{aligned}$$

The set \mathcal{G} of \mathfrak{S} generalizes to ψ any bounded uniformly continuous function.

Lemma 28.4.5 *If $\phi_{\mu_n}(t) \rightarrow \phi_{\mu}(t)$ where $\{\mu_n\}$ and μ are probability measures defined on the Borel sets of \mathbb{R}^p , then if ψ is any bounded uniformly continuous function,*

$$\lim_{n \rightarrow \infty} \int_{\mathbb{R}^p} \psi d\mu_n = \int_{\mathbb{R}^p} \psi d\mu.$$

Proof: Let $\varepsilon > 0$ be given, let ψ be a bounded function in $C^\infty(\mathbb{R}^p)$. Now let $\eta \in C_c^\infty(Q_r)$ where $Q_r \equiv [-r, r]^p$ satisfy the additional requirement that $\eta = 1$ on $Q_{r/2}$ and $\eta(x) \in [0, 1]$ for all x . By Lemma 28.4.3 the set, $\{\mu_n\}_{n=1}^\infty$, is tight and so if $\varepsilon > 0$ is given, there exists r sufficiently large such that for all n ,

$$\int_{[x \notin Q_{r/2}]} |1 - \eta| |\psi| d\mu_n < \frac{\varepsilon}{3},$$

and since μ is a single measure, the following holds whenever r is large enough.

$$\int_{[x \notin Q_{r/2}]} |1 - \eta| |\psi| d\mu < \frac{\varepsilon}{3}.$$

Thus,

$$\begin{aligned}\left| \int_{\mathbb{R}^p} \psi d\mu_n - \int_{\mathbb{R}^p} \psi d\mu \right| &\leq \left| \int_{\mathbb{R}^p} \psi d\mu_n - \int_{\mathbb{R}^p} \psi \eta d\mu_n \right| + \\ &\left| \int_{\mathbb{R}^p} \psi \eta d\mu_n - \int_{\mathbb{R}^p} \psi \eta d\mu \right| + \left| \int_{\mathbb{R}^p} \psi \eta d\mu - \int_{\mathbb{R}^p} \psi d\mu \right| \\ &\leq \frac{2\varepsilon}{3} + \left| \int_{\mathbb{R}^p} \psi \eta d\mu_n - \int_{\mathbb{R}^p} \psi \eta d\mu \right| < \varepsilon\end{aligned}$$

whenever n is large enough by Lemma 28.4.4 because $\psi\eta \in \mathfrak{S}$. This establishes the conclusion of the lemma in the case where ψ is also infinitely differentiable. To consider the general case, let ψ only be uniformly continuous and let $\psi_k = \psi * \phi_k$ where ϕ_k is a mollifier whose support is in $(-(1/k), (1/k))^p$. Then ψ_k converges uniformly to ψ and so the desired conclusion follows for ψ after a routine estimate. \blacksquare

Here are some items which are of considerable interest for their own sake.

Theorem 28.4.6 *Let H be a compact metric space. Then there exists a compact subset of $[0, 1]$, K and a continuous function, θ which maps K onto H .*

Proof: Without loss of generality, it can be assumed H is an infinite set since otherwise the conclusion is trivial. You could pick finitely many points of $[0, 1]$ for K .

Since H is compact, it is totally bounded. Therefore, there exists a 1 net for H $\{h_i\}_{i=1}^{m_1}$. Letting $H_i^1 \equiv \overline{B(h_i, 1)}$, it follows H_i^1 is also a compact metric space and so there exists a 1/2 net for each H_i^1 , $\{h_j^i\}_{j=1}^{m_i}$. Then taking the intersection of $\overline{B(h_j^i, \frac{1}{2})}$ with H_i^1 to obtain sets

denoted by H_j^2 and continuing this way, one can obtain compact subsets of H , $\{H_k^i\}$ which satisfies: each H_j^i is contained in some H_k^{i-1} , each H_j^i is compact with diameter less than i^{-1} , each H_j^i is the union of sets of the form H_k^{i+1} which are contained in it. Denoting by $\{H_j^i\}_{j=1}^{m_i}$ those sets corresponding to a superscript of i , it can also be assumed $m_i < m_{i+1}$.

If this is not so, simply add in another point to the i^{-1} net. Now let $\{I_j^i\}_{j=1}^{m_i}$ be disjoint closed intervals in $[0, 1]$ each of length no longer than 2^{-m_i} which have the property that I_j^i is contained in I_k^{i-1} for some k . Letting $K_i \equiv \bigcup_{j=1}^{m_i} I_j^i$, it follows K_i is a sequence of nested compact sets. Let $K = \bigcap_{i=1}^{\infty} K_i$. Then each $x \in K$ is the intersection of a unique sequence of these closed intervals, $\{I_{j_k}^k\}_{k=1}^{\infty}$. Define $\theta x \equiv \bigcap_{k=1}^{\infty} H_{j_k}^k$. Since the diameters of the H_j^i converge to 0 as $i \rightarrow \infty$, this function is well defined. It is continuous because if $x_n \rightarrow x$, then ultimately x_n and x are both in $I_{j_k}^k$, the k^{th} closed interval in the sequence whose intersection is x . Hence,

$$d(\theta x_n, \theta x) \leq \text{diameter}(H_{j_k}^k) \leq 1/k.$$

To see the map is onto, let $h \in H$. Then from the construction, there exists a sequence $\{H_{j_k}^k\}_{k=1}^{\infty}$ of the above sets whose intersection equals h . Then $\theta \left(\bigcap_{i=1}^{\infty} I_{j_k}^k \right) = h$. ■

Note θ is probably not one to one.

As an important corollary, it follows that the continuous functions defined on any compact metric space is separable.

Corollary 28.4.7 *Let H be a compact metric space and let $C(H)$ denote the continuous functions defined on H with the usual norm,*

$$\|f\|_{\infty} \equiv \max \{|f(x)| : x \in H\}$$

Then $C(H)$ is separable.

Proof: The proof is by contradiction. Suppose $C(H)$ is not separable. Let \mathcal{H}_k denote a maximal collection of functions of $C(H)$ with the property that if $f, g \in \mathcal{H}_k$, then $\|f - g\|_{\infty} \geq 1/k$. The existence of such a maximal collection of functions is a consequence of a simple use of the Hausdorff maximality theorem. Then $\bigcup_{k=1}^{\infty} \mathcal{H}_k$ is dense. Therefore, it cannot be countable by the assumption that $C(H)$ is not separable. It follows that for some k , \mathcal{H}_k is uncountable. Now by Theorem 28.4.6 there exists a continuous function θ defined on a compact subset K of $[0, 1]$ which maps K onto H . Now consider the functions defined on K

$$\mathcal{G}_k \equiv \{f \circ \theta : f \in \mathcal{H}_k\}.$$

Then \mathcal{G}_k is an uncountable set of continuous functions defined on K with the property that the distance between any two of them is at least as large as $1/k$. This contradicts separability of $C(K)$ which follows from the Weierstrass approximation theorem in which the separable countable set of functions is the restrictions of polynomials that involve only rational coefficients. ■

The next theorem gives the existence of a measure based on an assumption that a set of measures is tight. It is a sort of sequential compactness result. It is Prokhorov's theorem about a tight set of measures. Recall that Λ is tight means that for every $\varepsilon > 0$ there exists K compact such that $\mu(K^c) < \varepsilon$ for all $\mu \in \Lambda$.

Theorem 28.4.8 Let $\Lambda = \{\mu_n\}_{n=1}^\infty$ be a sequence of probability measures defined on the Borel sets of E a Polish space. If Λ is tight then there exists a probability measure λ and a subsequence of $\{\mu_n\}_{n=1}^\infty$, still denoted by $\{\mu_n\}_{n=1}^\infty$ such that whenever ϕ is a continuous bounded complex valued function defined on E ,

$$\lim_{n \rightarrow \infty} \int \phi d\mu_n = \int \phi d\lambda.$$

Conversely, if μ_n converges weakly to λ , then $\{\mu_n\}$ is tight.

Proof: By tightness, there exists an increasing sequence of compact sets, $\{K_n\}$ such that $\mu(K_n) > 1 - \frac{1}{n}$ for all $\mu \in \Lambda$. Now letting $\mu \in \Lambda$ and $\phi \in C(K_n)$ such that $\|\phi\|_\infty \leq 1$, it follows

$$\left| \int_{K_n} \phi d\mu \right| \leq \mu(K_n) \leq 1$$

and so the restrictions of the measures of Λ to K_n are contained in the unit ball of $C(K_n)'$. Recall from the Riesz representation theorem, the dual space of $C(K_n)$ is a space of complex Borel measures. Theorem 21.5.5 on Page 557 implies the unit ball of $C(K_n)'$ is weak * sequentially compact. This follows from the observation that $C(K_n)$ is separable which follows easily from the Weierstrass approximation theorem. Recall from the Riesz representation theorem, the dual space of $C(K_n)$ is a space of complex Borel measures. Theorem 21.5.5 on Page 557 implies the unit ball of $C(K_n)'$ is weak * sequentially compact. This follows from the observation that $C(K_n)$ is separable which is proved in Corollary 28.4.7 and leads to the fact that the unit ball in $C(K_n)'$ is actually metrizable by Theorem 21.5.5 on Page 557.

Thus the unit ball in $C(K_n)'$ is actually metrizable by Theorem 21.5.5 on Page 557. Therefore, there exists a subsequence of Λ , $\{\mu_{1k}\}$ such that their restrictions to K_1 converge weak * to a measure, $\lambda_1 \in C(K_1)'$. That is, for every $\phi \in C(K_1)$,

$$\lim_{k \rightarrow \infty} \int_{K_1} \phi d\mu_{1k} = \int_{K_1} \phi d\lambda_1$$

By the same reasoning, there exists a further subsequence $\{\mu_{2k}\}$ such that the restrictions of these measures to K_2 converge weak * to a measure $\lambda_2 \in C(K_2)'$ etc. Continuing this way,

$$\begin{aligned} \mu_{11}, \mu_{12}, \mu_{13}, \dots &\rightarrow \text{Weak * in } C(K_1)' \\ \mu_{21}, \mu_{22}, \mu_{23}, \dots &\rightarrow \text{Weak * in } C(K_2)' \\ \mu_{31}, \mu_{32}, \mu_{33}, \dots &\rightarrow \text{Weak * in } C(K_3)' \\ &\vdots \end{aligned}$$

Here the j^{th} sequence is a subsequence of the $(j-1)^{\text{th}}$. Let λ_n denote the measure in $C(K_n)'$ to which the sequence $\{\mu_{nk}\}_{k=1}^\infty$ converges weak *. Let $\{\mu_n\} \equiv \{\mu_{mn}\}$, the diagonal sequence. Thus this sequence is ultimately a subsequence of every one of the above sequences and so μ_n converges weak * in $C(K_m)'$ to λ_m for each m .

Claim: For $p > n$, the restriction of λ_p to the Borel sets of K_n equals λ_n .

Proof of claim: Let H be a compact subset of K_n . Then there are sets, V_l open in K_n which are decreasing and whose intersection equals H . This follows because this is a metric

space. Then let $H \prec \phi_l \prec V_l$. It follows

$$\begin{aligned}\lambda_n(V_l) &\geq \int_{K_n} \phi_l d\lambda_n = \lim_{k \rightarrow \infty} \int_{K_n} \phi_l d\mu_k \\ &= \lim_{k \rightarrow \infty} \int_{K_p} \phi_l d\mu_k = \int_{K_p} \phi_l d\lambda_p \geq \lambda_p(H).\end{aligned}$$

Now considering the ends of this inequality, let $l \rightarrow \infty$ and pass to the limit to conclude $\lambda_n(H) \geq \lambda_p(H)$. Similarly,

$$\begin{aligned}\lambda_n(H) &\leq \int_{K_n} \phi_l d\lambda_n = \lim_{k \rightarrow \infty} \int_{K_n} \phi_l d\mu_k \\ &= \lim_{k \rightarrow \infty} \int_{K_p} \phi_l d\mu_k = \int_{K_p} \phi_l d\lambda_p \leq \lambda_p(V_l).\end{aligned}$$

Then passing to the limit as $l \rightarrow \infty$, it follows $\lambda_n(H) \leq \lambda_p(H)$. Thus the restriction of $\lambda_p, \lambda_p|_{K_n}$ to the compact sets of K_n equals λ_n . Then by inner regularity it follows the two measures, $\lambda_p|_{K_n}$, and λ_n are equal on all Borel sets of K_n . Recall that for finite measures on the Borel sets of separable metric spaces, regularity is obtained for free.

It is fairly routine to exploit regularity of the measures to verify that $\lambda_m(F) \geq 0$ for all F a Borel subset of K_m . (Whenever $\phi \geq 0$, $\int_{K_m} \phi d\lambda_m \geq 0$ because $\int_{K_m} \phi d\mu_k \geq 0$. Now you can approximate \mathcal{X}_F with a suitable nonnegative ϕ using regularity of the measure.) Also, letting $\phi \equiv 1$,

$$1 \geq \lambda_m(K_m) \geq 1 - \frac{1}{m}. \quad (28.9)$$

Define for F a Borel set,

$$\lambda(F) \equiv \lim_{n \rightarrow \infty} \lambda_n(F \cap K_n).$$

The limit exists because the sequence on the right is increasing due to the above observation that $\lambda_n = \lambda_m$ on the Borel subsets of K_m whenever $n > m$. Thus for $n > m$

$$\lambda_n(F \cap K_n) \geq \lambda_n(F \cap K_m) = \lambda_m(F \cap K_m).$$

Now let $\{F_k\}$ be a sequence of disjoint Borel sets. Then

$$\begin{aligned}\lambda(\cup_{k=1}^{\infty} F_k) &\equiv \lim_{n \rightarrow \infty} \lambda_n(\cup_{k=1}^{\infty} F_k \cap K_n) = \lim_{n \rightarrow \infty} \lambda_n(\cup_{k=1}^{\infty} (F_k \cap K_n)) \\ &= \lim_{n \rightarrow \infty} \sum_{k=1}^{\infty} \lambda_n(F_k \cap K_n) = \sum_{k=1}^{\infty} \lambda(F_k)\end{aligned}$$

the last equation holding by the monotone convergence theorem.

It remains to verify $\lim_{k \rightarrow \infty} \int \phi d\mu_k = \int \phi d\lambda$ for every ϕ bounded and continuous. This is where tightness is used again. Suppose $\|\phi\|_{\infty} < M$. Then as noted above, $\lambda_n(K_n) = \lambda(K_n)$ because for $p > n$, $\lambda_p(K_n) = \lambda_n(K_n)$ and so letting $p \rightarrow \infty$, the above is obtained. Also, from 28.9,

$$\begin{aligned}\lambda(K_n^C) &= \lim_{p \rightarrow \infty} \lambda_p(K_n^C \cap K_p) \leq \limsup_{p \rightarrow \infty} (\lambda_p(K_p) - \lambda_p(K_n)) \\ &\leq \limsup_{p \rightarrow \infty} (\lambda_p(K_p) - \lambda_n(K_n)) \leq \limsup_{p \rightarrow \infty} \left(1 - \left(1 - \frac{1}{n}\right)\right) = \frac{1}{n}\end{aligned}$$

Consequently,

$$\begin{aligned}
 \left| \int \phi d\mu_k - \int \phi d\lambda \right| &\leq \left| \int_{K_n^C} \phi d\mu_k + \int_{K_n} \phi d\mu_k - \left(\int_{K_n} \phi d\lambda + \int_{K_n^C} \phi d\lambda \right) \right| \\
 &\leq \left| \int_{K_n} \phi d\mu_k - \int_{K_n} \phi d\lambda_n \right| + \left| \int_{K_n^C} \phi d\mu_k - \int_{K_n^C} \phi d\lambda \right| \\
 &\leq \left| \int_{K_n} \phi d\mu_k - \int_{K_n} \phi d\lambda_n \right| + \left| \int_{K_n^C} \phi d\mu_k \right| + \left| \int_{K_n^C} \phi d\lambda \right| \\
 &\leq \left| \int_{K_n} \phi d\mu_k - \int_{K_n} \phi d\lambda_n \right| + \frac{M}{n} + \frac{M}{n}
 \end{aligned}$$

First let n be so large that $2M/n < \varepsilon/2$ and then pick k large enough that the above expression is less than ε .

Now suppose μ_n converges to λ weakly. Then for ε there is a compact set such that $\lambda(K) > 1 - \varepsilon/2$. This is true because of Lemma 9.8.5 on Page 255 which says that finite measures on a Polish space are inner regular. Then let ψ be a continuous function with values in $[0, 1]$ which equals 1 on K and is 0 off a compact set $\tilde{K} \supseteq K$. Then $\int \psi d\lambda > 1 - \varepsilon/2$ and also, there exists N such that for all $n \geq N$, $\int \psi d\mu_n > 1 - \varepsilon/2$. Thus $n \geq N$ implies $\mu_n(\tilde{K}) > 1 - \varepsilon/2$. Therefore, enlarging \tilde{K} finitely many times, one obtains $\tilde{K} \supseteq K$ such that for all μ_n and λ , $\lambda(\tilde{K}), \mu_n(\tilde{K}) > 1 - \varepsilon/2$. Thus $\mu_n(\tilde{K}^C) \leq \varepsilon/2 < \varepsilon$ for all n and so $\{\mu_n\}$ is tight as claimed. ■

Definition 28.4.9 Let $\mu, \{\mu_n\}$ be probability measures defined on the Borel sets of \mathbb{R}^p and let the sequence of probability measures, $\{\mu_n\}$ satisfy

$$\lim_{n \rightarrow \infty} \int \phi d\mu_n = \int \phi d\mu.$$

for every ϕ a bounded continuous function. Then μ_n is said to converge weakly to μ .

With the above, it is possible to prove the following amazing theorem of Levy.

Theorem 28.4.10 Suppose $\{\mu_n\}$ is a sequence of probability measures defined on the Borel sets of \mathbb{R}^p and let $\{\phi_{\mu_n}\}$ denote the corresponding sequence of characteristic functions. If there exists ψ which is continuous at $\mathbf{0}$, $\psi(\mathbf{0}) = 1$, and for all \mathbf{t} ,

$$\phi_{\mu_n}(\mathbf{t}) \rightarrow \psi(\mathbf{t}),$$

then there exists a probability measure λ defined on the Borel sets of \mathbb{R}^p and

$$\phi_\lambda(\mathbf{t}) = \psi(\mathbf{t}).$$

That is, ψ is a characteristic function of a probability measure. Also, $\{\mu_n\}$ converges weakly to λ .

Proof: By Lemma 28.4.3 $\{\mu_n\}$ is tight. Therefore, there exists a subsequence $\{\mu_{n_k}\}$ converging weakly to a probability measure λ which implies that

$$\phi_\lambda(\mathbf{t}) \equiv \int e^{it \cdot \mathbf{x}} d\lambda(\mathbf{x}) = \lim_{n \rightarrow \infty} \int e^{it \cdot \mathbf{x}} d\mu_{n_k}(\mathbf{x}) = \lim_{n \rightarrow \infty} \phi_{\mu_{n_k}}(\mathbf{t}) = \psi(\mathbf{t}) \quad \blacksquare$$

Note how it was only necessary to assume $\psi(\mathbf{0}) = 1$ and ψ is continuous at $\mathbf{0}$ in order to conclude that ψ is a characteristic function. This helps to see why Prokhorov's and Levy's theorems are so amazing. Limits of characteristic functions tend to be characteristic functions. What about random variables?

If you have a probability measure λ on the Borel sets of \mathbb{R}^p , is there a random variable \mathbf{X} such that $\lambda = \lambda_{\mathbf{X}}$? Yes. You could let $\Omega = \mathbb{R}^p$ and $X(\mathbf{x}) = \mathbf{x}$ and $P(E) \equiv \lambda(E)$ for all E Borel. Then $\lambda_{\mathbf{X}}(\mathbf{X}^{-1}(E)) \equiv P(E) \equiv \lambda(E)$ so this is indeed a random variable such that $\lambda = \lambda_{\mathbf{X}}$. Thus for a probability measure λ , you can generally get a random variable which has λ as its distribution measure. Later, this is considered more. You might have more than one random variable having λ as its distribution measure.

In this next corollary, it suffices to have the random variables have values in a Banach space. However, I will write $|\mathbf{X}|$ rather than $\|\mathbf{X}\|$.

Corollary 28.4.11 *In the context of Theorem 28.4.10, suppose μ_n is the distribution measure of the random variable X_n and that $\sup_n E(|X_n|^q) = M_q < \infty$ for all $q \geq 1$ and that μ_n converges weakly to the probability measure μ . Then if μ is the distribution measure for a random variable \mathbf{X} , then $E(|\mathbf{X}|^q) < \infty$ for all $q \geq 1$.*

Proof:

$$\begin{aligned} E(|\mathbf{X}|^q) &= \int_0^\infty P(|\mathbf{X}|^q > \alpha) d\alpha = \int_0^\infty \mu(|x|^q > \alpha) d\alpha \\ &\leq \int_0^\infty \mu(|x|^q > \alpha) d\alpha \leq \int_0^\infty \int_{\mathbb{R}^p} (1 - \psi_\alpha) d\mu d\alpha \end{aligned}$$

where $\psi_\alpha = 1$ on $B(\mathbf{0}, \frac{1}{2}\alpha^{1/q})$ is nonnegative, and is in $C_c(B(\mathbf{0}, \alpha^{1/q}))$. Thus if $|x|^q > \alpha$, then $1 - \psi_\alpha(x) = 1$ which shows the above inequality holds. Also, if $(1 - \psi_\alpha(x)) > 0$, then $|x| > \frac{1}{2}\alpha^{1/q}$ and so $|x|^q > \frac{1}{2^q}\alpha$. Since weak convergence holds and $1 - \psi_\alpha$ is a bounded continuous function,

$$\int_{\mathbb{R}^p} (1 - \psi_\alpha) d\mu = \lim_{n \rightarrow \infty} \int_{\mathbb{R}^p} (1 - \psi_\alpha) d\mu_n \quad (28.10)$$

Therefore, from the above and Fatou's lemma,

$$\begin{aligned} \int_{\Omega} |\mathbf{X}|^q dP &\leq \int_0^\infty \lim_{n \rightarrow \infty} \int_{\mathbb{R}^p} (1 - \psi_\alpha) d\mu_n d\alpha \\ &\leq \liminf_{n \rightarrow \infty} \int_0^\infty \mu_n \left(\left[|x|^q > \frac{1}{2^q} \alpha \right] \right) d\alpha \end{aligned}$$

Changing the variable,

$$= \liminf_{n \rightarrow \infty} 2^q \int_0^\infty \mu_n(|x|^q > \delta) d\delta = \liminf_{n \rightarrow \infty} 2^q E(|X_n|^q) < \infty \blacksquare$$

Now recall the multivariate normal distribution.

Definition 28.4.12 *A random vector \mathbf{X} , with values in \mathbb{R}^p has a multivariate normal distribution written as*

$$\mathbf{X} \sim N_p(\mathbf{m}, \Sigma)$$

if for all Borel $E \subseteq \mathbb{R}^p$, the distribution measure is given by

$$\lambda_{\mathbf{X}}(E) = \int_{\mathbb{R}^p} \mathcal{X}_E(\mathbf{x}) \frac{1}{(2\pi)^{p/2} \det(\Sigma)^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\mathbf{m})^* \Sigma^{-1}(\mathbf{x}-\mathbf{m})} d\mathbf{x}$$

for \mathbf{m} a given vector and Σ a given positive definite symmetric matrix. Recall also that the characteristic function of this random variable is

$$E(e^{it \cdot \mathbf{X}}) = e^{it \cdot \mathbf{m}} e^{-\frac{1}{2} t^* \Sigma t} \quad (28.11)$$

So what if $\det(\Sigma) = 0$? Is there a probability measure having characteristic function $e^{it \cdot \mathbf{m}} e^{-\frac{1}{2} t^* \Sigma t}$? Let $\Sigma_n \rightarrow \Sigma$ in the Frobenius norm, $\det(\Sigma_n) > 0$. That is the $i j^{th}$ components converge and all the eigenvalues are positive. Then from the definition of the characteristic function,

$$\phi_{\lambda_{\mathbf{X}_n}}(\mathbf{t}) = e^{it \cdot \mathbf{m}} e^{-\frac{1}{2} t^* \Sigma_n t} \rightarrow \psi(\mathbf{t}) \equiv e^{it \cdot \mathbf{m}} e^{-\frac{1}{2} t^* \Sigma t}$$

Now clearly $\psi(\mathbf{0}) = 1$ and ψ is continuous so by Levy's theorem, Theorem 28.4.10, there is a probability measure μ such that $\psi(\mathbf{t}) = \phi_\mu(\mathbf{t})$. As noted above, there is also a random variable \mathbf{X} with $\lambda_{\mathbf{X}} = \mu$. Consider the moments for \mathbf{X}_n .

Lemma 28.4.13 *Let $\mathbf{X} \sim N(\mathbf{0}, \Sigma)$ where Σ is positive definite. Then the moments of \mathbf{X} all exist and are dominated by an expression which is continuously dependent on $\det(\Sigma)$.*

Proof: Let $q \geq 1$. $E(|\mathbf{X}|^q) = \int_{\mathbb{R}^p} \frac{|\mathbf{x}|^q}{(2\pi)^{p/2} \det(\Sigma)^{1/2}} e^{-\frac{1}{2} \mathbf{x}^* \Sigma^{-1} \mathbf{x}} d\mathbf{x}$. Let R be an orthogonal matrix with $\Sigma = R^* D R$ where D is a diagonal matrix having the positive eigenvalues σ_j on the diagonal. Thus $\mathbf{x}^* \Sigma^{-1} \mathbf{x} = \mathbf{x}^* R^* D^{-1} R \mathbf{x}$ so let $R \mathbf{x} \equiv \mathbf{y}$. Changing the variable in the integral and assuming $q = 2m$ for m a positive integer,

$$\begin{aligned} E(|\mathbf{X}|^{2m}) &= \int_{\mathbb{R}^p} \frac{(\sum_{k=1}^p y_k^2)^m}{(2\pi)^{p/2} \prod_{j=1}^p \sigma_j^{1/2}} e^{-\frac{1}{2} \mathbf{y}^* D^{-1} \mathbf{y}} d\mathbf{y} \\ &= 2^p \frac{1}{(2\pi)^{p/2}} \int_0^\infty \cdots \int_0^\infty \frac{(\sum_{k=1}^p y_k^2)^m}{\prod_{j=1}^p \sigma_j^{1/2}} e^{-\frac{1}{2} \mathbf{y}^* D^{-1} \mathbf{y}} d\mathbf{y} \end{aligned}$$

From convexity of $x \rightarrow x^m$

$$\begin{aligned} &= 2^p \frac{1}{(2\pi)^{p/2}} \int_0^\infty \cdots \int_0^\infty \frac{(p \sum_{k=1}^p \frac{1}{p} y_k^2)^m}{\prod_{j=1}^p \sigma_j^{1/2}} e^{-\frac{1}{2} \mathbf{y}^* D^{-1} \mathbf{y}} d\mathbf{y} \\ &\leq 2^p \frac{p^{m-1}}{(2\pi)^{p/2}} \int_0^\infty \cdots \int_0^\infty \frac{\sum_{k=1}^p y_k^{2m}}{\prod_{j=1}^p \sigma_j^{1/2}} e^{-\frac{1}{2} \sum_{k=1}^p y_k^2 \sigma_k^{-1}} d\mathbf{y} \\ &= 2^p \frac{p^{m-1}}{(2\pi)^{p/2}} \int_0^\infty \cdots \int_0^\infty \frac{\sum_{k \neq l} y_k^{2m}}{\prod_{j=1}^p \sigma_j^{1/2}} e^{-\frac{1}{2} \sum_{k \neq l} y_k^2 \sigma_k^{-1}} dx_1 \cdots \widehat{dx_l} \cdots dx_p \\ &\quad \cdot \int_0^\infty \frac{1}{\sigma_l^{1/2}} y_l^{2m} e^{-\frac{1}{2} y_l^2 \sigma_l^{-1}} dy_l \end{aligned}$$

Now letting $u = y_l \sigma_l^{-1/2}$, $dy_l = \sigma_l^{1/2} du$ and so that last integral is of the form

$$\int_0^\infty \sigma_l^m u^{2m} e^{-\frac{1}{2} u^2} du = \hat{C}_m \sigma_l^m$$

and so, doing this repeatedly, one obtains for the above integral an expression of the form

$$2^p \frac{p^{m-1}}{(2\pi)^{p/2}} C_m \left(\prod_{k=1}^p \sigma_k \right)^m = 2^p \frac{p^{m-1}}{(2\pi)^{p/2}} C_m \det(\Sigma)^m$$

which shows that the moments of \mathbf{X} all exist and are dominated by an expression which depends continuously on $\det(\Sigma)$. ■

In particular, these moments are bounded in case $\Sigma_n \rightarrow \Sigma$ where perhaps $\det(\Sigma) = 0$ but Σ_n is positive definite. With Corollary 28.4.11, this has proved the following theorem about the generalized normal distribution.

Theorem 28.4.14 *Let Σ be nonnegative and self adjoint $p \times p$ matrix. Then there exists a random variable \mathbf{X} whose distribution measure $\lambda_{\mathbf{X}}$ has characteristic function $\psi(\mathbf{t}) \equiv e^{-\frac{1}{2}\mathbf{t}^*\Sigma\mathbf{t}}$. Then all the moments exist and $E(\mathbf{X}\mathbf{X}^*) = \Sigma$.*

Proof: It remains to verify $E(\mathbf{X}\mathbf{X}^*) = \Sigma$ but this is routine from the fact that the moments exist. Use the characteristic function to compute $E(X_i X_j)$. Take $\frac{d}{dt_j} \left(\frac{d}{dt_i} (\psi(\mathbf{t})) \right)$. Using repeated index summation convention,

$$\begin{aligned} \psi(\mathbf{t}) &= e^{-\frac{1}{2}tr\Sigma_{rs}t_s}, \psi_{t_i} = e^{-\frac{1}{2}tr\Sigma_{rs}t_s} (-\Sigma_{is}t_s), \psi_{t_i t_j} \\ &= e^{-\frac{1}{2}tr\Sigma_{rs}t_s} (-\Sigma_{js}t_s) (-\Sigma_{is}t_s) + e^{-\frac{1}{2}tr\Sigma_{rs}t_s} (-\Sigma_{ij}) \end{aligned}$$

Thus $i^2 E(X_i X_j) = -\Sigma_{ij}$ showing that $E(\mathbf{X}\mathbf{X}^*) = \Sigma$ as claimed. ■

The case where $\mathbf{m} = \mathbf{0}$ is the one of most interest here, but you could always reduce to this case by considering a random variable $\mathbf{X} - \mathbf{m}$ where $E(\mathbf{X}) = \mathbf{m}$. There is an interesting corollary to this theorem.

Corollary 28.4.15 *Let H be a real Hilbert space. Then there exist random variables $W(h)$ for $h \in H$ such that for any finite set $\{f_1, f_2, \dots, f_n\}$,*

$$(W(f_1), W(f_2), \dots, W(f_n))$$

is normally distributed with mean 0 and covariance $\Sigma_{ij} = (f_i, f_j)$ and for every h, g ,

$$E(W(h)W(g)) = (h, g)_H$$

If $\{e_i\}$ is an orthogonal set of vectors of H , then $\{W(e_i)\}$ are independent random variables.

Proof: Let $\mu_{h_1 \dots h_m}$ be a generalized multivariate normal probability distribution with covariance $\Sigma_{ij} = (h_i, h_j)$ and mean 0. That such a thing exists follows from Theorem 28.4.14. Thus the characteristic function of this probability measure is $e^{-\frac{1}{2}\mathbf{t}^*\Sigma\mathbf{t}}$. Now consider $E_{k_1} \times \dots \times E_{k_n}$ for Borel sets E_{k_j} where $\{h_1, \dots, h_m\} \subseteq \{k_1 \dots k_n\}$ for $n > m$ and the set $E_{k_j} = \mathbb{R}$ whenever $k_j \notin \{h_1, \dots, h_m\}$. For simplicity, say h_1, \dots, h_m are the first m slots of k_1, \dots, k_n . Now consider $\mu_{k_1 \dots k_n}$,

$$\{h_1 \dots h_m, k_{m+1} \dots k_n\} = \{k_1 \dots k_n\}$$

Let ν be a measure on $\mathcal{B}(\mathbb{R}^m)$ which is given by $\nu(E) \equiv \mu_{k_1 \dots k_n}(E \times \mathbb{R}^{n-m})$. Then does it follow that $\nu = \mu_{h_1 \dots h_m}$? If so, then the Kolmogorov consistency condition will hold for

these measures $\mu_{h_1 \cdots h_m}$. To determine whether this is so, take the characteristic function of ν . Let Σ_1 be the $n \times n$ matrix which comes from the $\{k_1 \cdots k_n\}$ and let Σ_2 be the one which comes from the $\{h_1 \cdots h_m\}$.

$$\begin{aligned} \int_{\mathbb{R}^m} e^{it \cdot x} d\nu(x) &\equiv \int_{\mathbb{R}^{n-m}} \int_{\mathbb{R}^m} e^{i(t, \mathbf{0}) \cdot (x, y)} d\mu_{k_1 \cdots k_n}(x, y) \\ &= e^{-\frac{1}{2}(t^*, \mathbf{0}^*)\Sigma_1(t, \mathbf{0})} = e^{-\frac{1}{2}t^*\Sigma_2 t} \end{aligned}$$

which is the characteristic function for $\mu_{h_1 \cdots h_m}$. Therefore, these two measures are the same and the Kolmogorov consistency condition holds. It follows from The Kolmogorov extension theorem Theorem 20.3.3 that there exists a measure μ defined on the Borel sets of $\prod_{h \in H} \mathbb{R}$ which extends all of these measures. This argument also shows that if a random vector \mathbf{X} has characteristic function $e^{-\frac{1}{2}t^*\Sigma t}$, then if X_k is one of its components, then the characteristic function of X_k is $e^{-\frac{1}{2}t^2|h_k|^2}$ so this scalar valued random variable has mean zero and variance $|h_k|^2$. Then if $\omega \in \prod_{h \in H} \mathbb{R}$, $W(h)(\omega) \equiv \pi_h(\omega)$ where π_h denotes the projection onto position h in this product. Also define

$$(W(f_1), W(f_2), \dots, W(f_n)) \equiv \pi_{f_1 \cdots f_n}(\omega)$$

Then this is a random variable whose covariance matrix is just $\Sigma_{ij} = (f_i, f_j)_H$ and whose characteristic equation is $e^{-\frac{1}{2}t^*\Sigma t}$ so this verifies that

$$(W(f_1), W(f_2), \dots, W(f_n))$$

is normally distributed with covariance Σ . If you have two of them, $W(g), W(h)$, then $E(W(h)W(g)) = (h, g)_H$. This follows from what was just shown that $(W(f), W(g))$ is normally distributed and so the covariance will be

$$\begin{pmatrix} |f|^2 & (f, g) \\ (f, g) & |g|^2 \end{pmatrix} = \begin{pmatrix} E(W(f)^2) & E(W(f)W(g)) \\ E(W(f)W(g)) & E(W(g)^2) \end{pmatrix}$$

Finally consider the claim about independence. Any finite subset of $\{W(e_i)\}$ is generalized normal with the covariance matrix being a diagonal. Therefore,

$$(W(e_{i_1}), \dots, W(e_{i_n}))$$

is normally distributed with covariance a diagonal matrix so by Theorem 28.2.3, the random variables $\{W(e_i)\}$ are independent. ■

28.5 The Central Limit Theorem

The central limit theorem is one of the most marvelous theorems in mathematics. It can be proved through the use of characteristic functions. Recall for $x \in \mathbb{R}^p$,

$$\|x\|_\infty \equiv \max \{|x_j|, j = 1, \dots, p\}.$$

Also recall the definition of the distribution function for a random vector, \mathbf{X} .

$$F_{\mathbf{X}}(x) \equiv P(X_j \leq x_j, j = 1, \dots, p).$$

How can you tell if a sequence of random vectors with values in \mathbb{R}^p is tight? The next lemma gives a way to do this. It is Lemma 28.4.3. I am stating it here for convenience.

Lemma 28.5.1 *If $\{X_n\}$ is a sequence of random vectors with values in \mathbb{R}^p such that*

$$\lim_{n \rightarrow \infty} \phi_{X_n}(t) \equiv \lim_{n \rightarrow \infty} \phi_{\lambda_{X_n}}(t) = \psi(t)$$

for all t , where $\psi(0) = 1$ and ψ is continuous at 0 , then $\{\lambda_{X_n}\}_{n=1}^{\infty}$ is tight.

In proving the central limit theorem, one considers the pointwise convergence of characteristic functions and then seeks to obtain information about the distribution of the limit function. In fact, one is in the situation of the following lemma which is Lemma 28.4.4.

Lemma 28.5.2 *If $\phi_{X_n}(t) \rightarrow \phi_X(t)$ for all t , then whenever $\psi \in \mathfrak{S}$,*

$$\lambda_{X_n}(\psi) \equiv \int_{\mathbb{R}^p} \psi(y) d\lambda_{X_n}(y) \rightarrow \int_{\mathbb{R}^p} \psi(y) d\lambda_X(y) \equiv \lambda_X(\psi)$$

as $n \rightarrow \infty$.

The above gives what I want for $\psi \in \mathfrak{S}$ but this needs to be generalized to ψ any bounded uniformly continuous function. The following is Lemma 28.4.5.

Lemma 28.5.3 *If $\phi_{X_n}(t) \rightarrow \phi_X(t)$, then if ψ is any bounded uniformly continuous function,*

$$\lim_{n \rightarrow \infty} \int_{\mathbb{R}^p} \psi d\lambda_{X_n} = \int_{\mathbb{R}^p} \psi d\lambda_X.$$

Definition 28.5.4 *Let μ be a Radon measure on \mathbb{R}^p . A Borel set A , is a μ continuity set if $\mu(\partial A) = 0$ where $\partial A \equiv \bar{A} \setminus \text{int}(A)$ and int denotes the interior.*

The main result is the following continuity theorem. More can be said about the equivalence of various criteria [6].

Theorem 28.5.5 *If $\phi_{X_n}(t) \rightarrow \phi_X(t)$ then $\lambda_{X_n}(A) \rightarrow \lambda_X(A)$ whenever A is a λ_X continuity set.*

Proof: First suppose K is a closed set and let

$$\psi_k(x) \equiv (1 - k \text{dist}(x, K))^+.$$

Thus, since K is closed $\lim_{k \rightarrow \infty} \psi_k(x) = \mathcal{X}_K(x)$. Choose k large enough that

$$\int_{\mathbb{R}^p} \psi_k d\lambda_X \leq \lambda_X(K) + \varepsilon.$$

Then by Lemma 28.5.3, applied to the bounded uniformly continuous function ψ_k ,

$$\limsup_{n \rightarrow \infty} \lambda_{X_n}(K) \leq \limsup_{n \rightarrow \infty} \int \psi_k d\lambda_{X_n} = \int \psi_k d\lambda_X \leq \lambda_X(K) + \varepsilon.$$

Since ε is arbitrary, this shows $\limsup_{n \rightarrow \infty} \lambda_{X_n}(K) \leq \lambda_X(K)$ for all K closed.

Next suppose V is open and let

$$\psi_k(x) = 1 - (1 - k \text{dist}(x, V^c))^+.$$

Thus $\psi_k(x) \in [0, 1]$, $\psi_k = 1$ if $\text{dist}(x, V^C) \geq 1/k$, and $\psi_k = 0$ on V^C . Since V is open, it follows $\lim_{k \rightarrow \infty} \psi_k(x) = \mathcal{X}_V(x)$. Choose k large enough that $\int \psi_k d\lambda_X \geq \lambda_X(V) - \varepsilon$. Then by Lemma 28.5.3,

$$\liminf_{n \rightarrow \infty} \lambda_{X_n}(V) \geq \liminf_{n \rightarrow \infty} \int \psi_k(x) d\lambda_{X_n} = \int \psi_k(x) d\lambda_X \geq \lambda_X(V) - \varepsilon$$

and since ε is arbitrary, $\liminf_{n \rightarrow \infty} \lambda_{X_n}(V) \geq \lambda_X(V)$. Now let $\lambda_X(\partial A) = 0$ for A a Borel set.

$$\begin{aligned} \lambda_X(\text{int}(A)) &\leq \liminf_{n \rightarrow \infty} \lambda_{X_n}(\text{int}(A)) \leq \liminf_{n \rightarrow \infty} \lambda_{X_n}(A) \leq \\ \limsup_{n \rightarrow \infty} \lambda_{X_n}(A) &\leq \limsup_{n \rightarrow \infty} \lambda_{X_n}(\bar{A}) \leq \lambda_X(\bar{A}). \end{aligned}$$

But $\lambda_X(\text{int}(A)) = \lambda_X(\bar{A})$ by assumption and so $\lim_{n \rightarrow \infty} \lambda_{X_n}(A) = \lambda_X(A)$ as claimed. ■

As an application of this theorem the following is a version of the central limit theorem in the situation in which the limit distribution is multivariate normal. It concerns a sequence of random vectors, $\{X_k\}_{k=1}^\infty$, which are identically distributed, have finite mean \mathbf{m} , and satisfy $E(|X_k|^2) < \infty$.

Definition 28.5.6 For X a random vector with values in \mathbb{R}^p , let

$$F_X(x) \equiv P(\{X_j \leq x_j \text{ for each } j = 1, 2, \dots, p\}).$$

A different proof of the central limit theorem is in [48].

Lemma 28.5.7 If all the z_i and w_i have absolute value no more than 1, then $|\prod_{i=1}^n z_i - \prod_{i=1}^n w_i| \leq \sum_{k=1}^n |z_k - w_k|$.

Proof: It is clearly true if $n = 1$. Suppose true for n . Then

$$\begin{aligned} \left| \prod_{i=1}^{n+1} z_i - \prod_{i=1}^{n+1} w_i \right| &\leq \left| \prod_{i=1}^{n+1} z_i - z_{n+1} \prod_{i=1}^n w_i \right| + \left| z_{n+1} \prod_{i=1}^n w_i - \prod_{i=1}^{n+1} w_i \right| \\ &\leq \left| \prod_{i=1}^n z_i - \prod_{i=1}^n w_i \right| + |z_{n+1} - w_{n+1}| \left| \prod_{i=1}^n w_i \right| \leq \sum_{k=1}^{n+1} |z_k - w_k| \quad \blacksquare \end{aligned}$$

Theorem 28.5.8 Let $\{X_k\}_{k=1}^\infty$ be random vectors satisfying $E(|X_k|^2) < \infty$ which are independent and identically distributed with mean \mathbf{m} and positive definite covariance $\Sigma \equiv E((X - \mathbf{m})(X - \mathbf{m})^*)$. Let $Z_n \equiv \sum_{j=1}^n \frac{X_j - \mathbf{m}}{\sqrt{n}}$. Then for $Z \sim N_p(\mathbf{0}, \Sigma)$, $\lim_{n \rightarrow \infty} F_{Z_n}(x) = F_Z(x)$ for all x .

Proof: The characteristic function of Z_n is given by

$$\phi_{Z_n}(t) = E\left(e^{it \cdot \sum_{j=1}^n \frac{X_j - \mathbf{m}}{\sqrt{n}}}\right) = \prod_{j=1}^n E\left(e^{it \cdot \left(\frac{X_j - \mathbf{m}}{\sqrt{n}}\right)}\right).$$

By Taylor's theorem applied to real and imaginary parts of e^{ix} , it follows

$$e^{ix} = 1 + ix - f(x) \frac{x^2}{2}$$

where $|f(x)| < 2$ and $\lim_{x \rightarrow 0} f(x) = 1$. Denoting \mathbf{X}_j as \mathbf{X} , this implies

$$e^{it \cdot \left(\frac{\mathbf{X} - \mathbf{m}}{\sqrt{n}} \right)} = 1 + it \cdot \frac{\mathbf{X} - \mathbf{m}}{\sqrt{n}} - f \left(t \cdot \left(\frac{\mathbf{X} - \mathbf{m}}{\sqrt{n}} \right) \right) \frac{(t \cdot (\mathbf{X} - \mathbf{m}))^2}{2n}$$

Thus $e^{it \cdot \left(\frac{\mathbf{X} - \mathbf{m}}{\sqrt{n}} \right)} = 1 + it \cdot \frac{\mathbf{X} - \mathbf{m}}{\sqrt{n}} - \frac{(t \cdot (\mathbf{X} - \mathbf{m}))^2}{2n} + \left(1 - f \left(t \cdot \left(\frac{\mathbf{X} - \mathbf{m}}{\sqrt{n}} \right) \right) \right) \frac{(t \cdot (\mathbf{X} - \mathbf{m}))^2}{2n}$. This implies

$$\phi_{\mathbf{Z}_n}(t) = \prod_{j=1}^n E \left[1 - \frac{(t \cdot (\mathbf{X}_j - \mathbf{m}))^2}{2n} + \frac{(t \cdot (\mathbf{X}_j - \mathbf{m}))^2}{2n} \left(1 - f \left(t \cdot \left(\frac{\mathbf{X}_j - \mathbf{m}}{\sqrt{n}} \right) \right) \right) \right]$$

Then $\phi_{\mathbf{Z}_n}(t) =$

$$\begin{aligned} & \prod_{j=1}^n E \left[1 - \frac{(t \cdot (\mathbf{X}_j - \mathbf{m}))^2}{2n} + \frac{(t \cdot (\mathbf{X}_j - \mathbf{m}))^2}{2n} \left(1 - f \left(t \cdot \left(\frac{\mathbf{X}_j - \mathbf{m}}{\sqrt{n}} \right) \right) \right) \right] \\ & - \prod_{j=1}^n E \left[1 - \frac{(t \cdot (\mathbf{X}_j - \mathbf{m}))^2}{2n} \right] + \prod_{j=1}^n \left(1 - \frac{E(t \cdot (\mathbf{X}_j - \mathbf{m}))^2}{2n} \right) \end{aligned}$$

Now $(t \cdot (\mathbf{X} - \mathbf{m}))^2 = \mathbf{t}^* (\mathbf{X} - \mathbf{m}) (\mathbf{X} - \mathbf{m})^* \mathbf{t}$. Since these \mathbf{X}_k are identically distributed with the same mean \mathbf{m} , the above is of the form

$$e_n + \prod_{j=1}^n \left(1 - \frac{E(t \cdot (\mathbf{X}_j - \mathbf{m}))^2}{2n} \right) = e_n + \left(1 - \frac{1}{2n} \mathbf{t}^* \Sigma \mathbf{t} \right)^n$$

where for large n , the needed expressions have small absolute value and so, from the above lemma, for large n ,

$$|e_n| \leq \frac{1}{2n} \sum_{j=1}^n E \left((t \cdot (\mathbf{X}_j - \mathbf{m}))^2 \left| 1 - f \left(t \cdot \left(\frac{\mathbf{X}_j - \mathbf{m}}{\sqrt{n}} \right) \right) \right| \right)$$

Now write \mathbf{X} for \mathbf{X}_k since all are identically distributed. Then the above right side is no more than

$$\frac{1}{2} E \left((t \cdot (\mathbf{X} - \mathbf{m}))^2 \left| 1 - f \left(t \cdot \left(\frac{\mathbf{X} - \mathbf{m}}{\sqrt{n}} \right) \right) \right| \right)$$

which converges to 0 as $n \rightarrow \infty$ by the dominated convergence theorem. Therefore,

$$\lim_{n \rightarrow \infty} \phi_{\mathbf{Z}_n}(t) = \lim_{n \rightarrow \infty} \left(1 - \frac{1}{2n} \mathbf{t}^* \Sigma \mathbf{t} \right)^n = e^{-\frac{1}{2} \mathbf{t}^* \Sigma \mathbf{t}} = \phi_{\mathbf{Z}}(t)$$

where $\mathbf{Z} \sim N_p(\mathbf{0}, \Sigma)$. Therefore, from Theorem 28.5.5, $F_{\mathbf{Z}_n}(x) \rightarrow F_{\mathbf{Z}}(x)$ for all x because

$$R_x \equiv \prod_{k=1}^p (-\infty, x_k]$$

is a set of λ_Z continuity due to the assumption that $\lambda_Z \ll m_p$ which is implied by $Z \sim N_p(\mathbf{0}, \Sigma)$. ■

Suppose \mathbf{X} is a random vector with covariance Σ and mean \mathbf{m} , and suppose also that Σ^{-1} exists. Consider $\Sigma^{-(1/2)}(\mathbf{X} - \mathbf{m}) \equiv \mathbf{Y}$. Then $E(\mathbf{Y}) = \mathbf{0}$ and

$$\begin{aligned} E(\mathbf{Y}\mathbf{Y}^*) &= E\left(\Sigma^{-(1/2)}(\mathbf{X} - \mathbf{m})(\mathbf{X}^* - \mathbf{m})\Sigma^{-(1/2)}\right) \\ &= \Sigma^{-(1/2)}E((\mathbf{X} - \mathbf{m})(\mathbf{X}^* - \mathbf{m}))\Sigma^{-(1/2)} = I. \end{aligned}$$

Thus \mathbf{Y} has zero mean and covariance I . This implies the following corollary to Theorem 28.5.8.

Corollary 28.5.9 *Let $\{X_j\}_{j=1}^\infty$ be independent identically distributed random variables and suppose they have mean \mathbf{m} and positive definite covariance Σ where Σ^{-1} exists. Then if*

$$\mathbf{Z}_n \equiv \sum_{j=1}^n \Sigma^{-(1/2)} \frac{(\mathbf{X}_j - \mathbf{m})}{\sqrt{n}},$$

it follows that for $\mathbf{Z} \sim N_p(\mathbf{0}, I)$, $F_{\mathbf{Z}_n}(\mathbf{x}) \rightarrow F_{\mathbf{Z}}(\mathbf{x})$ for all \mathbf{x} .

Chapter 29

Martingales

29.1 Conditional Expectation

From Observation 27.4.5 on Page 748, it was shown that the conditional expectation of a random variable X given some others really is just what the words suggest. Given $\omega \in \Omega$, it results in a value for the “other” random variables and then you essentially take the expectation of X given this information which yields the value of the conditional expectation of X given the other random variables. It was also shown in Lemma 27.4.4 that this gives the same result as finding a $\sigma(X_1, \dots, X_n)$ measurable function Z such that for all $F \in \sigma(X_1, \dots, X_n)$,

$$\int_F X dP = \int_F Z dP$$

This was done for a particular type of σ algebra but there is no need to be this specialized. The following is the general version of conditional expectation given a σ algebra. It makes perfect sense to ask for the conditional expectation given a σ algebra and this is what will be done from now on.

Definition 29.1.1 Let (Ω, \mathcal{M}, P) be a probability space and let $\mathcal{S} \subseteq \mathcal{F}$ be two σ algebras contained in \mathcal{M} . Let f be \mathcal{F} measurable and in $L^1(\Omega; W)$ for W a Banach space. Then $E(f|\mathcal{S})$, called the conditional expectation of f with respect to \mathcal{S} is defined as follows:

$$E(f|\mathcal{S}) \text{ is } \mathcal{S} \text{ measurable}$$

For all $E \in \mathcal{S}$,

$$\int_E E(f|\mathcal{S}) dP = \int_E f dP$$

The existence and uniqueness of the conditional expectation is described earlier in Theorem 24.12.1 on Page 702. For convenience, here is this theorem.

Theorem 29.1.2 Let E be a separable Banach space and $X \in L^1(\Omega; E, \mathcal{M})$ where X is measurable with respect to \mathcal{M} and let \mathcal{S} be a σ algebra which is contained in \mathcal{M} . Then there exists a unique $Z \in L^1(\Omega; E, \mathcal{S})$ such that for all $A \in \mathcal{S}$,

$$\int_A X dP = \int_A Z dP$$

Denoting this Z as $E(X|\mathcal{S})$, it follows

$$\|E(X|\mathcal{S})\| \leq E(\|X\| |\mathcal{S}).$$

A few properties are described next. Let W be a separable Banach space in the following lemma.

Lemma 29.1.3 The above is well defined. Also, if $\mathcal{S} \subseteq \mathcal{F}$ then if $X \in L^1(\Omega; W)$,

$$E(X|\mathcal{S}) = E(E(X|\mathcal{F})|\mathcal{S}). \quad (29.1)$$

If Z is in $L^\infty(\Omega; W')$ bounded and measurable in \mathcal{S} then

$$ZE(X|\mathcal{S}) = E(ZX|\mathcal{S}). \quad (29.2)$$

Also, if $a, b \in W'$, and $X, Y \in L^1(\Omega; W)$

$$aE(X|\mathcal{S}) + bE(Y|\mathcal{S}) = E(aX + bY|\mathcal{S}). \quad (29.3)$$

Proof: To begin with consider 29.3. By definition, if $F \in \mathcal{S}$,

$$\begin{aligned} \int_F aE(X|\mathcal{S}) + bE(Y|\mathcal{S}) dP &= a \int_F E(X|\mathcal{S}) dP + b \int_F E(Y|\mathcal{S}) dP \\ &= a \int_F X dP + b \int_F Y dP = \int_F (aX + bY) dP \equiv \int_F E(aX + bY|\mathcal{S}) dP \end{aligned}$$

Since F is arbitrary, this shows 29.3.

Let $F \in \mathcal{S}$. Then

$$\begin{aligned} \int_F E(E(X|\mathcal{F})|\mathcal{S}) dP &\equiv \int_F E(X|\mathcal{F}) dP \\ &\equiv \int_F X dP \equiv \int_F E(X|\mathcal{S}) dP \end{aligned}$$

and so, by uniqueness, $E(E(X|\mathcal{F})|\mathcal{S}) = E(X|\mathcal{S})$. This shows 29.1.

To establish 29.2, note that if $Z = a\mathcal{X}_F$ where $F \in \mathcal{S}$, and $a \in W'$, by Definition 29.1.1,

$$\begin{aligned} \int a\mathcal{X}_F E(X|\mathcal{S}) dP &= \int_F E(aX|\mathcal{S}) dP dP = \int_F aX dP \\ &= \int a\mathcal{X}_F X dP = \int E(a\mathcal{X}_F X|\mathcal{S}) dP \end{aligned}$$

which shows 29.2 in the case where Z is $a\mathcal{X}_F$, $F \in \mathcal{S}$. It follows this also holds for simple functions with values in W' . Let Z be in $L^\infty(\Omega; W)$. By Theorem 24.2.4 there is a sequence of simple functions $\{s_n\}$, $\|s_n(\omega)\| \leq 2\|Z(\omega)\|$ which converges to Z and let $F \in \mathcal{S}$. Then by what was just shown,

$$\int_F s_n E(X|\mathcal{S}) dP = \int_F E(s_n X|\mathcal{S}) dP \equiv \int_F s_n X dP \quad (29.4)$$

Now

$$\begin{aligned} \left\| \int_F E(s_n X|\mathcal{S}) dP - \int_F E(ZX|\mathcal{S}) dP \right\| &= \left\| \int_F (s_n - Z) X dP \right\| \\ &\leq \int_F \|(s_n - Z) X\| dP \end{aligned}$$

and this converges to 0 by the dominated convergence theorem. Also from Theorem 24.12.1

$$\|s_n E(X|\mathcal{S})\| = \|E(s_n X|\mathcal{S})\| \leq E(\|s_n X\||\mathcal{S}) \leq 2E(\|ZX\||\mathcal{S})$$

which is in $L^1(\Omega)$. Thus one can apply the dominated convergence theorem to the left side of 29.4 and use what was just shown to pass to a limit in 29.4 and obtain

$$\int_F ZE(X|\mathcal{S}) dP = \int_F ZX dP \equiv \int_F E(ZX|\mathcal{S}) dP.$$

Since this holds for every $F \in \mathcal{S}$, this shows 29.2. ■

The next major result is a generalization of Jensen's inequality whose proof depends on the following lemma about convex functions. It pertains to the case where the functions have values in \mathbb{R} .

Lemma 29.1.4 *Let ϕ be a convex real valued function defined on an interval I . Then for each $x \in I$, there exists a_x such that for all $t \in I$,*

$$\phi(t) \geq a_x(t-x) + \phi(x).$$

Also ϕ is continuous on I .

Proof: Let $x \in I$ and let $t > x$. Then by convexity of ϕ ,

$$\frac{\phi(x + \lambda(t-x)) - \phi(x)}{\lambda(t-x)} \leq \frac{\phi(x)(1-\lambda) + \lambda\phi(t) - \phi(x)}{\lambda(t-x)} = \frac{\phi(t) - \phi(x)}{t-x}.$$

Therefore $t \rightarrow \frac{\phi(t) - \phi(x)}{t-x}$ is increasing if $t > x$. If $t < x, t-x < 0$ so

$$\frac{\phi(x + \lambda(t-x)) - \phi(x)}{\lambda(t-x)} \geq \frac{\phi(x)(1-\lambda) + \lambda\phi(t) - \phi(x)}{\lambda(t-x)} = \frac{\phi(t) - \phi(x)}{t-x}$$

and so $t \rightarrow \frac{\phi(t) - \phi(x)}{t-x}$ is increasing for $t \neq x$. Let

$$a_x \equiv \inf \left\{ \frac{\phi(t) - \phi(x)}{t-x} : t > x \right\}.$$

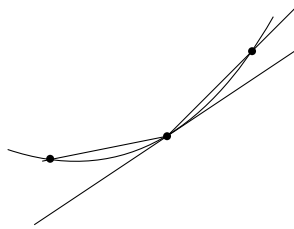
Then if $t_1 < x$, and $t > x$,

$$\frac{\phi(t_1) - \phi(x)}{t_1 - x} \leq a_x \leq \frac{\phi(t) - \phi(x)}{t-x}.$$

Thus for all $t \in I$,

$$\phi(t) \geq a_x(t-x) + \phi(x). \quad (29.5)$$

The continuity of ϕ follows easily from this and the observation that convexity simply says that the graph of ϕ lies below the line segment joining two points on its graph. Thus, we have the following picture which clearly implies continuity. ■



Lemma 29.1.5 *Let I be an interval on \mathbb{R} and let ϕ be a convex function defined on I . Then there exists a sequence $\{(a_n, b_n)\}$ such that*

$$\phi(t) = \sup \{a_n t + b_n, n = 1, \dots\}.$$

Proof: Let a_x be as defined in the above lemma. Let

$$\psi(x) \equiv \sup \{a_r(x-r) + \phi(r) : r \in \mathbb{Q} \cap I\}.$$

Thus if $r_1 \in \mathbb{Q}$, $\psi(r_1) \equiv \sup \{a_r(r_1 - r) + \phi(r) : r \in \mathbb{Q} \cap I\} \geq \phi(r_1)$. Then ψ is convex on I so ψ is continuous. Therefore, $\psi(t) \geq \phi(t)$ for all $t \in I$. By 29.5,

$$\psi(t) \geq \phi(t) \geq \sup \{a_r(t - r) + \phi(r), r \in \mathbb{Q} \cap I\} \equiv \psi(t).$$

Thus $\psi(t) = \phi(t)$. Let $\mathbb{Q} \cap I = \{r_n\}$, $a_n = a_{r_n}$ and $b_n = -a_{r_n}r_n + \phi(r_n)$. Continuity gives the desired results at endpoints of I . ■

In this lemma, X, Y have values in \mathbb{R} .

Lemma 29.1.6 *If $X \leq Y$, then $E(X|\mathcal{S}) \leq E(Y|\mathcal{S})$ a.e. Also $X \rightarrow E(X|\mathcal{S})$ is linear.*

Proof: Let $A \in \mathcal{S}$.

$$\int_A E(X|\mathcal{S}) dP \equiv \int_A X dP \leq \int_A Y dP \equiv \int_A E(Y|\mathcal{S}) dP.$$

Hence $E(X|\mathcal{S}) \leq E(Y|\mathcal{S})$ a.e. as claimed. That $X \rightarrow E(X|\mathcal{S})$ is linear follows from Lemma 29.1.3.

Theorem 29.1.7 *(Jensen's inequality) Let $X(\omega) \in I$ a closed interval and let $\phi : I \rightarrow \mathbb{R}$ be convex. Suppose $E(|X|), E(|\phi(X)|) < \infty$. Then $\phi(E(X|\mathcal{S})) \leq E(\phi(X)|\mathcal{S})$.*

Proof: Let $\phi(x) = \sup \{a_n x + b_n\}$. Letting $A \in \mathcal{S}$,

$$\frac{1}{P(A)} \int_A E(X|\mathcal{S}) dP = \frac{1}{P(A)} \int_A X dP \in I \text{ a.e.}$$

whenever $P(A) \neq 0$. The claim that $\frac{1}{P(A)} \int_A X dP \in I$ follows from approximating X with simple functions having values in I . Hence $E(X|\mathcal{S})(\omega) \in I$ a.e. and so it makes sense to consider $\phi(E(X|\mathcal{S}))$. Now $a_n E(X|\mathcal{S}) + b_n = E(a_n X + b_n|\mathcal{S}) \leq E(\phi(X)|\mathcal{S})$. Thus $\sup \{a_n E(X|\mathcal{S}) + b_n\} = \phi(E(X|\mathcal{S})) \leq E(\phi(X)|\mathcal{S})$ a.e. ■

29.2 Conditional Expectation and Independence

The situation of interest is a sequence of random variables $\{Y_i\}$ having values in a separable Banach space along with two other random variables X, Z both of which are measurable with respect to \mathcal{E} where $\mathcal{E}, \sigma(Y_1, \dots)$ are independent σ algebras contained in \mathcal{M} , so if $A \in \mathcal{E}$ and $B \in \sigma(Y_1, \dots)$, then $P(A \cap B) = P(A)P(B)$. Also let $\sigma(Z, Y)$ be the smallest σ algebra for which Z and each Y_k are measurable. Then both X, Z would seem to relate only to \mathcal{E} and so it would seem that the values of the Y_k would be irrelevant and $E(X|Z) = E(X|\sigma(Z, Y))$. That is, the conditional expectation given the extra conditions from the Y_k is unchanged. Is it like the earlier notion in which independence means you can dispense with the givens? Recall also the notation $E(X|Z)$ is defined as $E(X|\sigma(Z))$.

Recall that if $f^{-1}(O) \in \mathcal{F}$ a σ algebra and this holds for all O open, then if $\mathcal{S} \equiv \{B : f^{-1}(B) \in \mathcal{F}\}$ it follows that \mathcal{S} is a σ algebra and so it contains the Borel sets.

Proposition 29.2.1 *Let $\mathcal{E}, \sigma(Y_1, Y_2, \dots)$ be independent σ algebras contained in \mathcal{M} . Also let $\sigma(Z, Y)$ be the smallest σ algebra which respect to which each Y_k and Z is measurable. Let Z, X be \mathcal{E} measurable. Then $E(X|Z) = E(X|\sigma(Z, Y))$.*

Proof: First I claim that $\sigma(Z, Y)$ consists of $(Z, Y)^{-1}(B)$ where B is Borel in $W \times W^{\mathbb{N}}$ where $W^{\mathbb{N}} \equiv \prod_{i=1}^{\infty} W$. To see this, let

$$\left\{ (Z, Y)^{-1}(B) : B \text{ is Borel in } W^{1+\mathbb{N}} \right\} \equiv \mathcal{F}.$$

Then \mathcal{F} is a σ algebra and for V_J of the form $\prod_{j=1}^{\infty} V_j$ where $V_j = W$ for all j except for those contained in a finite set J and for the other j , V_j is an open set, and U an open set in W , then $(Z, Y)^{-1}(U \times V_J) = Z^{-1}(U) \cap Y^{-1}(V_J) \in \mathcal{F}$ so that, in particular, choosing V_J and U appropriately shows Z, Y_k are all measurable with respect to \mathcal{F} . Hence $\mathcal{F} \supseteq \sigma(Z, Y)$. Also, $U \times V_J$ just described, where U is in a countable basis for W and each V_j is W or in a countable basis for W is a countable basis for the topology of $W^{1+\mathbb{N}}$. By definition, $\sigma(Z, Y)$ must contain $(Z, Y)^{-1}(U \times V_J)$ for $U \times V_J$ in this countable basis. Therefore, $\sigma(Z, Y)$ must contain $(Z, Y)^{-1}(O)$ for all O open and so also $\sigma(Z, Y)$ must contain $(Z, Y)^{-1}(B)$ for B Borel, so $\sigma(Z, Y) \supseteq \mathcal{F}$. Also, this shows that for \mathcal{H} these sets in the countable basis for $W^{1+\mathbb{N}}$, $\sigma(\mathcal{H}) = \sigma(Z, Y) = \mathcal{F}$. Now consider the following computation.

$$\begin{aligned} & \int_{(Z, Y)^{-1}(U \times V_J)} E(X | \sigma(Z, Y)) dP \\ & \equiv \int_{(Z, Y)^{-1}(U \times V_J)} X dP \\ & = \int \mathcal{X}_{Y^{-1}(V_J)} \mathcal{X}_{Z^{-1}(U)} X dP \stackrel{*}{=} P(Y^{-1}(V_J)) \int \mathcal{X}_{Z^{-1}(U)} X dP \\ & \stackrel{**}{=} P(Y^{-1}(V_J)) \int E(\mathcal{X}_{Z^{-1}(U)} X | Z) dP = \int \mathcal{X}_{Y^{-1}(V_J)} E(\mathcal{X}_{Z^{-1}(U)} X | Z) dP \\ & = \int_{Z^{-1}(U)} \mathcal{X}_{Y^{-1}(V_J)} E(X | Z) dP = \int_{(Z, Y)^{-1}(U \times V_J)} E(X | Z) dP \end{aligned}$$

Then $*$ happens because $\mathcal{X}_{[Z \in U]} X$ is \mathcal{E} measurable and $\mathcal{X}_{Y^{-1}(V_J)}$ is \mathcal{F} measurable and by assumption, these are independent σ algebras. $**$ happens because Ω is $\sigma(Z)$ measurable and the definition of conditional expectation. Next step happens because $\mathcal{X}_{Y^{-1}(V_J)}$ is \mathcal{F} measurable and $E(\mathcal{X}_{Z^{-1}(U)} X | Z)$ is \mathcal{E} measurable. Then the rest follows because $\mathcal{X}_{Z^{-1}(U)}$ is $\sigma(Z)$ measurable so it comes out of the conditional expectation.

Now let \mathcal{G} be those sets G of $\sigma(Z, Y) = \sigma(\mathcal{H})$ such that

$$\int_G X dP \equiv \int_G E(X | \sigma(Z, Y)) dP = \int_G E(X | \sigma(Z)) dP \equiv \int_G E(X | Z) dP$$

As just shown, $\mathcal{H} \subseteq \mathcal{G}$. If $G = \cup_{k=1}^{\infty} G_k$ where the G_k are disjoint and the equation holds for each G_k

$$\begin{aligned} \int_G E(X | Z) dP &= \int \sum_{k=1}^{\infty} \mathcal{X}_{G_k} E(X | Z) dP = \sum_{k=1}^{\infty} \int_{G_k} E(X | Z) dP \\ &= \sum_{k=1}^{\infty} \int_{G_k} E(X | \sigma(Z, Y)) dP = \int \sum_{k=1}^{\infty} \mathcal{X}_{G_k} E(X | \sigma(Z, Y)) dP \\ &= \int_G E(X | \sigma(Z, Y)) dP \end{aligned}$$

Thus \mathcal{G} is closed with respect to countable disjoint unions. If $G \in \mathcal{G}$ then both sides of the following equal $E(X)$ and so

$$\begin{aligned} & \int_G E(X|\sigma(Z, Y)) dP + \int_{G^c} E(X|\sigma(Z, Y)) dP \\ &= \int_G E(X|Z) dP + \int_{G^c} E(X|Z) dP \end{aligned}$$

and subtracting $\int_G E(X|\sigma(Z, Y)) dP = \int_G E(X|Z) dP$ from both sides yields

$$\int_{G^c} E(X|\sigma(Z, Y)) dP = \int_{G^c} E(X|Z) dP.$$

This proves \mathcal{G} is all of $\sigma(\mathcal{K}) = \sigma(Z, Y)$. By uniqueness, $E(X|Z) = E(X|\sigma(Z, Y))$ a.e. ■

29.3 Discrete Stochastic Processes

Earlier a special case of a discrete martingale and sub-martingale was discussed. This section considers the general case where one just has an increasing list of σ algebras. The idea is that you have an increasing list of real numbers $\{a_n\}$ which is well ordered and $X(a_n) \equiv X_{a_n}$ is measurable with respect to \mathcal{F}_{a_n} where the \mathcal{F}_{a_n} are increasing in n . Such a sequence of \mathcal{F}_{a_n} measurable functions is called a stochastic processes. We usually let this well ordered increasing list be some subset of \mathbb{N} for the sake of convenience, but it could be equally spaced points in an interval, for example. We let \mathcal{F} denote a σ algebra which contains all of these \mathcal{F}_k . For convenience, you could call it \mathcal{F}_∞ .

Definition 29.3.1 Let \mathcal{F}_k be an increasing sequence of σ algebras which are subsets of \mathcal{F} and X_k be a sequence of Banach space valued random variables with $E(\|X_k\|) < \infty$ such that X_k is \mathcal{F}_k measurable. Such a thing is called a stochastic process. It is called a martingale if

$$E(X_{k+1}|\mathcal{F}_k) = X_k,$$

In case the Banach space is \mathbb{R} the stochastic process is a sub-martingale if

$$E(X_{k+1}|\mathcal{F}_k) \geq X_k,$$

and a supermartingale if

$$E(X_{k+1}|\mathcal{F}_k) \leq X_k.$$

Saying that X_k is \mathcal{F}_k measurable is referred to by saying $\{X_k\}$ is adapted to \mathcal{F}_k .

For sub and super martingales, you need to be considering X which has values in \mathbb{R} . No such restriction is necessary for a martingale. If $\{X_k\}$ is a Banach space valued martingale, then from Theorem 29.1.2, $\{\|X_k\|\}$ is a sub-martingale and that if $\{X_k\}$ is a real sub-martingale and ϕ is convex and **increasing**, then $\{\phi(X_k)\}$ is a sub-martingale. This is discussed below.

Also in general, for a stochastic process, $E(X_n|\mathcal{F}_{n-2}) = E(E(X_n|\mathcal{F}_{n-1})|\mathcal{F}_{n-2})$. Thus if $\{X_n\}$ is a sub-martingale, Lemma 29.1.6 implies $E(X_n|\mathcal{F}_{n-2}) \geq X_{n-2}$. Similarly,

$$E(X_n|\mathcal{F}_k) \geq X_k$$

whenever $k < n$. Something similar happens with a martingale where you can replace \geq with $=$ or a super-martingale where you replace \geq with \leq . I will use this observation without comment in what follows.

Lemma 29.3.2 *Let ϕ be a convex and increasing function and suppose*

$$\{(X_n, \mathcal{F}_n)\}$$

is a sub-martingale. Then if $E(|\phi(X_n)|) < \infty$, it follows

$$\{(\phi(X_n), \mathcal{F}_n)\}$$

is also a sub-martingale.

Proof: It is given that $E(X_{n+1}, \mathcal{F}_n) \geq X_n$ and so

$$\phi(X_n) \leq \phi(E(X_{n+1} | \mathcal{F}_n)) \leq E(\phi(X_{n+1}) | \mathcal{F}_n)$$

by Jensen's inequality. ■

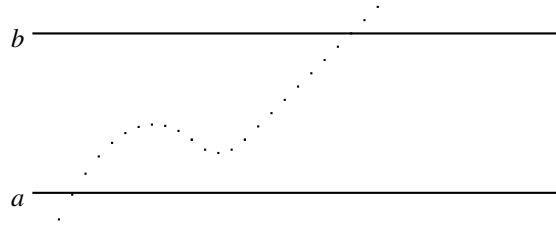
Certainly one of the most amazing things about sub-martingales is the convergence theorem. Recall the earlier sub-martingale convergence theorem. I am going to present the same thing here in this more general setting. Then later, I will present it again in the case of continuous sub-martingales. This which follows is almost identical to the earlier proof.

So why did I even bother with the earlier development? It is because I wanted to make a smooth transition from the idea we usually have of conditional probability where it is probability of something given the value of something else to this more general and much more abstract notion involving σ algebras.

An upcrossing occurs when a sequence goes from a up to b . Thus it crosses the interval, $[a, b]$ in the up direction, hence the name upcrossing. More precisely,

Definition 29.3.3 *Let $\{x_i\}_{i=1}^I$ be any sequence of real numbers, $I \leq \infty$. Define an increasing sequence of integers $\{m_k\}$ as follows. m_1 is the first integer ≥ 1 such that $x_{m_1} \leq a$, m_2 is the first integer larger than m_1 such that $x_{m_2} \geq b$, m_3 is the first integer larger than m_2 such that $x_{m_3} \leq a$, etc. Then each sequence, $\{x_{m_{2k-1}}, \dots, x_{m_{2k}}\}$, is called an upcrossing of $[a, b]$.*

Here is a picture of an upcrossing.



Proposition 29.3.4 *Let $\{X_i\}_{i=1}^n$ be a finite sequence of real random variables defined on Ω where (Ω, \mathcal{F}, P) is a probability space. Let $U_{[a,b]}(\omega)$ denote the number of upcrossings of $X_i(\omega)$ of the interval $[a, b]$. Then $U_{[a,b]}$ is a random variable.*

Proof: Let $X_0(\omega) \equiv a + 1$, let $Y_0(\omega) \equiv 0$, and let $Y_k(\omega)$ remain 0 for $k = 0, \dots, l$ until $X_l(\omega) \leq a$. When this happens (if ever), $Y_{l+1}(\omega) \equiv 1$. Then let $Y_i(\omega)$ remain 1 for $i = l + 1, \dots, r$ until $X_r(\omega) \geq b$ when $Y_{r+1}(\omega) \equiv 0$. Let $Y_k(\omega)$ remain 0 for $k \geq r + 1$ until $X_k(\omega) \leq a$ when $Y_k(\omega) \equiv 1$ and continue in this way. Thus the upcrossings of $X_i(\omega)$

are identified as unbroken strings of ones for Y_k with a zero at each end, with the possible exception of the last string of ones which may be missing the zero at the upper end and may or may not be an upcrossing.

Note also that Y_0 is measurable because it is identically equal to 0 and that if Y_k is measurable, then Y_{k+1} is measurable because the only change in going from k to $k+1$ is a change from 0 to 1 or from 1 to 0 on a measurable set determined by X_k . In particular,

$$Y_{k+1}^{-1}(1) = ([Y_k = 1] \cap [X_k < b]) \cup ([Y_k = 0] \cap [X_k \leq a])$$

This set is in \mathcal{F} by induction. Of course, $Y_{k+1}^{-1}(0)$ is just the complement of this set. Thus Y_{k+1} is \mathcal{F} measurable since 0, 1 are the only two values. Now let

$$Z_k(\omega) = \begin{cases} 1 & \text{if } Y_k(\omega) = 1 \text{ and } Y_{k+1}(\omega) = 0, \\ 0 & \text{otherwise,} \end{cases}$$

if $k < n$ and

$$Z_n(\omega) = \begin{cases} 1 & \text{if } Y_n(\omega) = 1 \text{ and } X_n(\omega) \geq b, \\ 0 & \text{otherwise.} \end{cases}$$

Thus $Z_k(\omega) = 1$ exactly when an upcrossing has been completed and each Z_i is a random variable.

$$U_{[a,b]}(\omega) = \sum_{k=1}^n Z_k(\omega)$$

so $U_{[a,b]}$ is a random variable as claimed. ■

The following corollary collects some key observations found in the above construction.

Corollary 29.3.5 $U_{[a,b]}(\omega) \leq$ the number of unbroken strings of ones in the sequence $\{Y_k(\omega)\}$, there being at most one unbroken string of ones which produces no upcrossing. Also

$$Y_i(\omega) = \psi_i\left(\{X_j(\omega)\}_{j=1}^{i-1}\right), \quad (29.6)$$

where ψ_i is some function of the past values of $X_j(\omega)$.

The following is called the upcrossing lemma.

29.3.1 Upcrossings

Lemma 29.3.6 (upcrossing lemma) Let $\{(X_i, \mathcal{F}_i)\}_{i=1}^n$ be a sub-martingale and let

$$U_{[a,b]}(\omega)$$

be the number of upcrossings of $[a, b]$. Then

$$E(U_{[a,b]}) \leq \frac{E(|X_n|) + |a|}{b - a}.$$

Proof: Let $\phi(x) \equiv a + (x - a)^+$ so that ϕ is an increasing convex function always at least as large as a . By Lemma 29.3.2 it follows that $\{(\phi(X_k), \mathcal{F}_k)\}$ is also a sub-martingale.

$$\phi(X_{k+r}) - \phi(X_k) = \sum_{i=k+1}^{k+r} \phi(X_i) - \phi(X_{i-1})$$

$$= \sum_{i=k+1}^{k+r} (\phi(X_i) - \phi(X_{i-1})) Y_i + \sum_{i=k+1}^{k+r} (\phi(X_i) - \phi(X_{i-1})) (1 - Y_i).$$

Observe that Y_i is \mathcal{F}_{i-1} measurable from its construction in Proposition 29.3.4, Y_i depending only on X_j for $j < i$.

Now let the unbroken strings of ones for $\{Y_i(\omega)\}$ be

$$\{k_1, \dots, k_1 + r_1\}, \{k_2, \dots, k_2 + r_2\}, \dots, \{k_m, \dots, k_m + r_m\} \quad (29.7)$$

where $m = V(\omega) \equiv$ the number of unbroken strings of ones in the sequence $\{Y_i(\omega)\}$. By Corollary 29.3.5 $V(\omega) \geq U_{[a,b]}(\omega)$.

$$\begin{aligned} & \phi(X_n(\omega)) - \phi(X_1(\omega)) \\ &= \sum_{k=1}^n (\phi(X_k(\omega)) - \phi(X_{k-1}(\omega))) Y_k(\omega) \\ & \quad + \sum_{k=1}^n (\phi(X_k(\omega)) - \phi(X_{k-1}(\omega))) (1 - Y_k(\omega)). \end{aligned}$$

The first sum in the above reduces to summing over the unbroken strings of ones because the terms in which $Y_i(\omega) = 0$ contribute nothing. Therefore, observing that for $x > a$, $\phi(x) = x$,

$$\begin{aligned} \phi(X_n(\omega)) - \phi(X_1(\omega)) &\geq U_{[a,b]}(\omega)(b-a) + 0 + \\ & \quad \sum_{k=1}^n (\phi(X_k(\omega)) - \phi(X_{k-1}(\omega))) (1 - Y_k(\omega)) \end{aligned} \quad (29.8)$$

where the zero on the right side results from a string of ones which does not produce an upcrossing. It is here that it is important that $\phi(x) \geq a$. Such a string begins with $\phi(X_k(\omega)) = a$ and results in an expression of the form $\phi(X_{k+m}(\omega)) - \phi(X_k(\omega)) \geq 0$ since $\phi(X_{k+m}(\omega)) \geq a$. If X_k had not been replaced with $\phi(X_k)$, it would have been possible for $\phi(X_{k+m}(\omega))$ to be less than a and the zero in the above could have been a negative number. This would have been inconvenient.

Next take the expectation of both sides in 29.8. This results in

$$\begin{aligned} E(\phi(X_n) - \phi(X_1)) &\geq (b-a)E(U_{[a,b]}) \\ & \quad + E\left(\sum_{k=1}^n (\phi(X_k) - \phi(X_{k-1}))(1 - Y_k)\right) \\ &\geq (b-a)E(U_{[a,b]}) \end{aligned}$$

The reason for the last inequality where the term at the end was dropped is

$$\begin{aligned} & E((\phi(X_k) - \phi(X_{k-1}))(1 - Y_k)) \\ &= E(E((\phi(X_k) - \phi(X_{k-1}))(1 - Y_k) | \mathcal{F}_{k-1})) \\ &= E((1 - Y_k)E(\phi(X_k) | \mathcal{F}_{k-1}) - (1 - Y_k)E(\phi(X_{k-1}) | \mathcal{F}_{k-1})) \\ &\geq E((1 - Y_k)(\phi(X_{k-1}) - \phi(X_{k-1}))) = 0. \end{aligned}$$

Recall that Y_k is \mathcal{F}_{k-1} measurable and that $(\phi(X_k), \mathcal{F}_k)$ is a sub-martingale. ■

The reason for this lemma is to prove the amazing sub-martingale convergence theorem.

29.3.2 The Sub-martingale Convergence Theorem

Theorem 29.3.7 (*sub-martingale convergence theorem*) Let

$$\{(X_i, \mathcal{F}_i)\}_{i=1}^{\infty}$$

be a sub-martingale with $K \equiv \sup E(|X_n|) < \infty$. Then there exists a random variable X , such that $E(|X|) \leq K$ and

$$\lim_{n \rightarrow \infty} X_n(\omega) = X(\omega) \text{ a.e.}$$

Proof: Let $a, b \in \mathbb{Q}$ and let $a < b$. Let $U_{[a,b]}^n(\omega)$ be the number of upcrossings of $\{X_i(\omega)\}_{i=1}^n$. Then let

$$U_{[a,b]}(\omega) \equiv \lim_{n \rightarrow \infty} U_{[a,b]}^n(\omega) = \text{number of upcrossings of } \{X_i\}.$$

By the upcrossing lemma,

$$E(U_{[a,b]}^n) \leq \frac{E(|X_n|) + |a|}{b-a} \leq \frac{K + |a|}{b-a}$$

and so by the monotone convergence theorem,

$$E(U_{[a,b]}) \leq \frac{K + |a|}{b-a} < \infty$$

which shows $U_{[a,b]}(\omega)$ is finite a.e., for all $\omega \notin S_{[a,b]}$ where $P(S_{[a,b]}) = 0$. Define

$$S \equiv \cup \{S_{[a,b]} : a, b \in \mathbb{Q}, a < b\}.$$

Then $P(S) = 0$ and if $\omega \notin S$, $\{X_k\}_{k=1}^{\infty}$ has only finitely many upcrossings of every interval having rational endpoints. For such ω it cannot be the case that

$$\limsup_{k \rightarrow \infty} X_k(\omega) > \liminf_{k \rightarrow \infty} X_k(\omega)$$

because then you could pick rational a, b such that $[a, b]$ is between the limsup and the liminf and there would be infinitely many upcrossings of $[a, b]$. Thus, for $\omega \notin S$,

$$\limsup_{k \rightarrow \infty} X_k(\omega) = \liminf_{k \rightarrow \infty} X_k(\omega) = \lim_{k \rightarrow \infty} X_k(\omega) \equiv X_{\infty}(\omega).$$

Letting $X_{\infty}(\omega) \equiv 0$ for $\omega \in S$, Fatou's lemma implies

$$\int_{\Omega} |X_{\infty}| dP = \int_{\Omega} \liminf_{n \rightarrow \infty} |X_n| dP \leq \liminf_{n \rightarrow \infty} \int_{\Omega} |X_n| dP \leq K \blacksquare$$

As a simple application, here is an easy proof of a nice theorem about convergence of sums of independent random variables.

Theorem 29.3.8 Let $\{X_k\}$ be a sequence of independent real valued random variables such that $E(|X_k|) < \infty$, $E(X_k) = 0$, and

$$\sum_{k=1}^{\infty} E(X_k^2) < \infty.$$

Then $\sum_{k=1}^{\infty} X_k$ converges a.e.

Proof: Let $\mathcal{F}_n \equiv \sigma(X_1, \dots, X_n)$. Consider $S_n \equiv \sum_{k=1}^n X_k$.

$$E(S_{n+1} | \mathcal{F}_n) = S_n + E(X_{n+1} | \mathcal{F}_n).$$

Letting $A \in \mathcal{F}_n$ it follows from independence that

$$\begin{aligned} \int_A E(X_{n+1} | \mathcal{F}_n) dP &\equiv \int_A X_{n+1} dP = \int_{\Omega} \mathcal{X}_A X_{n+1} dP \\ &= P(A) \int_{\Omega} X_{n+1} dP = 0 \end{aligned}$$

and so $E(S_{n+1} | \mathcal{F}_n) = S_n$. Therefore, $\{S_n\}$ is a martingale. Now using independence again,

$$E(|S_n|) \leq E(|S_n|^2) = \sum_{k=1}^n E(X_k^2) \leq \sum_{k=1}^{\infty} E(X_k^2) < \infty$$

and so $\{S_n\}$ is an L^1 bounded martingale. Therefore, it converges a.e. ■

Corollary 29.3.9 Let $\{X_k\}$ be a sequence of independent real valued random variables such that $E(|X_k|) < \infty$, $E(X_k) = m_k$, and

$$\sum_{k=1}^{\infty} E(|X_k - m_k|^2) < \infty.$$

Then $\sum_{k=1}^{\infty} (X_k - m_k)$ converges a.e.

This can be extended to the case where the random variables have values in a separable Hilbert space. Recall that for $\{e_k\}$ an orthonormal basis and $\sum_{k=1}^{\infty} |a_k|_H^2 < \infty$, $\sum_{k=1}^{\infty} a_k e_k \in H$ and the infinite sum makes sense. Also, for $x \in H$, $x = \sum_k (x, e_k) e_k$ the convergence in H .

Theorem 29.3.10 Let $\{X_k\}$ be a sequence of independent H valued random variables where H is a real separable Hilbert space such that $E(|X_k|_H) < \infty$, $E(X_k) = 0$, and $\sum_{k=1}^{\infty} E(|X_k|_H^2) < \infty$. Then $\sum_{k=1}^{\infty} X_k$ converges a.e.

Proof: Let $\{e_k\}$ be an orthonormal basis for H . Then $\{(X_n, e_k)_H\}_{n=1}^{\infty}$ are real valued, independent, and their mean equals 0. Also

$$\sum_{n=1}^{\infty} E(|(X_n, e_k)_H|^2) \leq \sum_{n=1}^{\infty} E(|X_n|_H^2) < \infty$$

and so from Theorem 29.3.8, the series, $\sum_{n=1}^{\infty} (X_n, e_k)_H$ converges a.e. Therefore, there exists a set of measure zero such that for ω not in this set, $\sum_n (X_n(\omega), e_k)_H$ converges for each k . For ω not in this exceptional set, define

$$Y_k(\omega) \equiv \sum_{n=1}^{\infty} (X_n(\omega), e_k)_H$$

Next define $S(\omega) \equiv \sum_{k=1}^{\infty} Y_k(\omega) e_k$. Of course it is not clear this even makes sense. I need to show $\sum_{k=1}^{\infty} |Y_k(\omega)|^2 < \infty$. Using the independence of the X_n

$$\begin{aligned} E(|Y_k|^2) &\leq \liminf_{N \rightarrow \infty} E\left(\left(\sum_{n=1}^N \sum_{m=1}^N (X_n, e_k)_H (X_m, e_k)_H\right)\right) \\ &= \liminf_{N \rightarrow \infty} E\left(\sum_{n=1}^N (X_n, e_k)_H^2\right) = \sum_{n=1}^{\infty} E((X_n, e_k)_H^2) \end{aligned}$$

the last from the monotone convergence theorem. Hence from the above,

$$\begin{aligned} E\left(\sum_k |Y_k|^2\right) &= \sum_k E(|Y_k|^2) \leq \sum_k \sum_n E\left((X_n, e_k)_H^2\right) \\ &= \sum_n E\left(\sum_k (X_n, e_k)^2\right) = \sum_n E(|X_n|_H^2) < \infty \end{aligned}$$

the last by assumption. Therefore, for ω off a set of measure zero, and for $Y_k(\omega) \equiv \sum_{n=1}^{\infty} (X_n(\omega), e_k)_H$ which exists a.e. by Theorem 29.3.8, $\sum_k |Y_k(\omega)|^2 < \infty$ and so for a.e. ω , $S(\omega) \equiv \sum_{k=1}^{\infty} Y_k(\omega) e_k$ makes sense. Thus for these ω

$$\begin{aligned} S(\omega) &= \sum_l (S(\omega), e_l) e_l = \sum_l Y_l(\omega) e_l \equiv \sum_l \sum_n (X_n(\omega), e_l)_H e_l \\ &= \sum_n \sum_l (X_n(\omega), e_l) e_l = \sum_n X_n(\omega). \blacksquare \end{aligned}$$

Now with this theorem, here is a strong law of large numbers.

Theorem 29.3.11 Suppose $\{X_k\}$ are independent random variables and

$$E(|X_k|) < \infty$$

for each k and $E(X_k) = m_k$. Suppose also

$$\sum_{j=1}^{\infty} \frac{1}{j^2} E(|X_j - m_j|^2) < \infty. \quad (29.9)$$

Then $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n (X_j - m_j) = 0$ a.e.

Proof: Consider the sum

$$\sum_{j=1}^{\infty} \frac{X_j - m_j}{j}.$$

This sum converges a.e. because of 29.9 and Theorem 29.3.10 applied to the random vectors $\left\{ \frac{X_j - m_j}{j} \right\}$. Therefore, from Lemma 26.8.4 it follows that for a.e. ω ,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n (X_j(\omega) - m_j) = 0. \blacksquare$$

The next corollary is often called the strong law of large numbers. It follows immediately from the above theorem.

Corollary 29.3.12 Suppose $\{X_j\}_{j=1}^{\infty}$ are independent random vectors, $\lambda_{X_i} = \lambda_{X_j}$ for all i, j having mean m and variance equal to

$$\sigma^2 \equiv \int_{\Omega} |X_j - m|^2 dP < \infty.$$

Then for a.e. $\omega \in \Omega$, $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n X_j(\omega) = m$

29.3.3 Doob Sub-martingale Estimates

Another very interesting result about sub-martingales is the Doob sub-martingale estimate. First is a technical lemma which is frequently useful in situations where we want to interchange order of integration. I shall likely use this lemma without comment occasionally.

Lemma 29.3.13 *If f is \mathcal{F} measurable and nonnegative then*

$$(\lambda, \omega) \rightarrow \mathcal{X}_{[f > \lambda]} \text{ is } \mathcal{F} \times \mathcal{B}(\mathbb{R}) \text{ measurable.}$$

Proof: Let s be a nonnegative simple function, $s(\omega) = \sum_{k=1}^n c_k \mathcal{X}_{E_k}(\omega)$ where we can let the sum be written such that the c_k are strictly increasing in k and these are the positive values of s . Also let $F_k = \cup_{i=k}^n E_i$.

$$\mathcal{X}_{[s > \lambda]} = \sum_{k=1}^n \mathcal{X}_{[c_{k-1}, c_k)}(\lambda) \mathcal{X}_{F_k}(\omega), \quad c_0 \equiv 0.$$

which is clearly product measurable. To see that this formula is valid, first consider the case where $\lambda \in [0, c_1)$. Then $\mathcal{X}_{[s > \lambda]} = 1$ on F_1 and 0 off F_1 . The first term of the right side equals 1 and the others are 0 due to $\mathcal{X}_{[c_{k-1}, c_k)}(\lambda)$. Thus the formula holds for such λ . Now suppose $\lambda \in [c_{j-1}, c_j)$. Then left side is 1 when $s(\omega) = c_l$ for some $l \geq c_j$. In this case, the right side has exactly one term equal to 1 and it is $\mathcal{X}_{[c_{j-1}, c_j)}(\lambda) \mathcal{X}_{F_j}(\omega)$. The remaining case is that $\lambda \geq c_n$. In this case, the right side equals 0 and the left side also equals 0 because $s(\omega)$ is never strictly larger than c_n . ■

For arbitrary $f \geq 0$ and measurable, there is an increasing sequence of simple functions s_n converging pointwise to f . Therefore,

$$\lim_{n \rightarrow \infty} \mathcal{X}_{[s_n > \lambda]} = \mathcal{X}_{[f > \lambda]}$$

and so $\mathcal{X}_{[f > \lambda]}$ is product measurable. ■

Theorem 29.3.14 *Let $\{(X_i, \mathcal{F}_i)\}_{i=1}^\infty$ be a sub-martingale. Then for $\lambda > 0$,*

$$P\left(\left[\max_{1 \leq k \leq n} X_k > \lambda\right]\right) \leq \frac{1}{\lambda} \int_{\Omega} \mathcal{X}_{[\max_{1 \leq k \leq n} X_k > \lambda]} X_n^+ dP \leq \frac{1}{\lambda} \int_{\Omega} X_n^+ dP$$

Proof: Let

$$\begin{aligned} A_1 &\equiv [X_1 > \lambda], A_2 \equiv [X_2 > \lambda] \setminus A_1, \\ \dots, A_k &\equiv [X_k > \lambda] \setminus \left(\cup_{i=1}^{k-1} A_i\right) \dots \end{aligned}$$

Thus each A_k is \mathcal{F}_k measurable, the A_k are disjoint, and their union equals

$$\left[\max_{1 \leq k \leq n} X_k > \lambda\right].$$

Therefore from the definition of a sub-martingale and Jensen's inequality,

$$P\left(\left[\max_{1 \leq k \leq n} X_k > \lambda\right]\right) = \sum_{k=1}^n P(A_k) \leq \frac{1}{\lambda} \sum_{k=1}^n \int_{A_k} X_k dP$$

$$\begin{aligned}
&\leq \frac{1}{\lambda} \sum_{k=1}^n \int_{A_k} E(X_n | \mathcal{F}_k) dP \leq \frac{1}{\lambda} \sum_{k=1}^n \int_{A_k} E(X_n | \mathcal{F}_k)^+ dP \\
&\leq \frac{1}{\lambda} \sum_{k=1}^n \int_{A_k} E(X_n^+ | \mathcal{F}_k) dP = \frac{1}{\lambda} \sum_{k=1}^n \int_{A_k} X_n^+ dP \leq \frac{1}{\lambda} \int_{\Omega} \mathcal{X}_{[\max_{1 \leq k \leq n} X_k > \lambda]} X_n^+ dP. \blacksquare
\end{aligned}$$

Now suppose X_k is a martingale with values in W a Banach space. Then from the theorem on conditional expectation, $E(\|X_{k+1}\| | \mathcal{F}_k) \geq \|E(X_{k+1} | \mathcal{F}_k)\| = \|X_k\|$. Thus $k \rightarrow \|X_k\|$ is a sub-martingale and so one gets the following interesting corollary.

Corollary 29.3.15 *Let X_n be a martingale with values in a Banach space W . Then for $\lambda > 0$,*

$$P\left(\left[\max_{1 \leq k \leq n} \|X_k\| > \lambda\right]\right) \leq \frac{1}{\lambda} \int_{\Omega} \mathcal{X}_{[\max_{1 \leq k \leq n} \|X_k\| > \lambda]} \|X_n\| dP \leq \frac{1}{\lambda} \int_{\Omega} \|X_n\| dP$$

Now suppose X_k is a martingale with values in W a Banach space. For $p > 1, k \rightarrow \|X_k\|^p$ is a sub-martingale because

$$E(\|X_{k+1}\|^p | \mathcal{F}_k) \geq (E(\|X_{k+1}\| | \mathcal{F}_k))^p \geq \|E(X_{k+1} | \mathcal{F}_k)\|^p = \|X_k\|^p$$

Therefore, from the definition of the Lebesgue integral of a positive function,

$$\int_{\Omega} \left(\max_{1 \leq k \leq n} \|X_k\|\right)^p dP = \int_{\Omega} \max_{1 \leq k \leq n} \|X_k\|^p dP = \int_0^{\infty} P\left(\left[\max_{1 \leq k \leq n} \|X_k\|^p > \lambda\right]\right) d\lambda$$

Change variables $\lambda = \mu^p$ and using the Doob estimate, Theorem 29.3.14,

$$= \int_0^{\infty} P\left(\left[\max_{1 \leq k \leq n} \|X_k\| > \lambda^{1/p}\right]\right) d\lambda = p \int_0^{\infty} P\left(\left[\max_{1 \leq k \leq n} \|X_k\| > \mu\right]\right) \mu^{p-1} d\mu$$

To save on notation, let $X_n^* \equiv \max_{1 \leq k \leq n} \|X_k\|$. Then using Lemma 29.3.13 as needed,

$$\leq \int_0^{\infty} \frac{p\mu^{p-1}}{\mu} \int_{\Omega} \mathcal{X}_{[X_n^* > \mu]} \|X_n\| dP d\mu = \int_{\Omega} \|X_n(\omega)\| \int_0^{\infty} \frac{p\mu^{p-1}}{\mu} \mathcal{X}_{[X_n^* > \mu]} d\mu dP$$

Then $p-1 = p/q$ where $1/p + 1/q = 1$,

$$\begin{aligned}
&\leq \int_{\Omega} \int_0^{X_n^*} p\mu^{p-2} \|X_n\| d\mu dP = \frac{p}{p-1} \int_{\Omega} (X_n^*)^{p/q} \|X_n\| dP \\
&\leq \frac{p}{p-1} \left(\int_{\Omega} (X_n^*)^p dP\right)^{1/q} \left(\int_{\Omega} \|X_n\|^p dP\right)^{1/p}
\end{aligned}$$

This proves the following version of the above Doob estimate.

Theorem 29.3.16 *Let $n \rightarrow X_n$ be a martingale with values in a Banach space W and $X_n \in L^p(\Omega; W)$, $p > 1$, then for $X_n^* \equiv \max_{k \leq n} \{ \|X_k\| \}$, then*

$$\int_{\Omega} (X_n^*)^p dP \leq \frac{p}{p-1} \left(\int_{\Omega} (X_n^*)^p dP\right)^{1/q} \left(\int_{\Omega} \|X_n\|^p dP\right)^{1/p}$$

In fact,

$$\left(\int_{\Omega} (X_n^*)^p dP\right)^{1/p} \leq \frac{p}{p-1} \left(\int_{\Omega} \|X_n\|^p dP\right)^{1/p}$$

Proof: The second claim is all which remains. Let n be given. Then let

$$A_1 \equiv \{\omega : \|X_1(\omega)\| \geq \|X_k(\omega)\| \text{ for } k \neq 1\}.$$

Let

$$A_2 \equiv \{\omega : \|X_2(\omega)\| \geq \|X_k(\omega)\|, k \neq 2\} \setminus A_1,$$

etc. Then these sets are disjoint and so

$$X_n^* = \sum_{k=1}^n \mathcal{X}_{A_k} \|X_k\|, \quad \int_{\Omega} (X_n^*)^p dP = \sum_{k=1}^n \int_{A_k} \|X_k\|^p dP < \infty$$

and so, we can divide both sides with $(\int_{\Omega} (X_n^*)^p)^{1/q}$ to obtain the last claim. ■

Later on I will consider the case of continuous sub-martingales and martingales. In this case, you have to work harder and one way is to use a stopping time. These stopping times are about to be discussed in the next section for discrete processes.

29.4 Optional Sampling and Stopping Times

I think that the optional sampling theorem of Doob is amazing. That is why it gets repeated quite a bit. It is one of those theorems that you read and when you get to the end, having followed the argument, you sit back and feel amazed at what you just went through. You ask yourself if it is really true or whether you made some mistake. At least this is how it affects me.

First it is necessary to define the notion of a stopping time. If you have an increasing sequence of σ algebras $\{\mathcal{F}_n\}$ and a process $\{X_n\}$ such that X_n is \mathcal{F}_n measurable, the idea of a stopping time τ is that τ is measurable and $X_{\min(\tau, n)}$ is a \mathcal{F}_n measurable function. In other words, by stopping with this stopping time, we preserve the \mathcal{F}_n measurability. It is customary to write $n \wedge \tau = \min(n, \tau)$. Thus, we want to have $X_{n \wedge \tau}^{-1}(O) \in \mathcal{F}_n$ where O is an open set in some metric space where X_n has its values and $a \wedge b$ means $\min(a, b)$.

$$X_{n \wedge \tau}^{-1}(O) = [\tau \leq n] \cap [\omega : X_{\tau(\omega)}(\omega) \in O] \cup [\tau > n] \cap [\omega : X_n(\omega) \in O]$$

Now

$$[\tau \leq n] \cap [\omega : X_{\tau(\omega)}(\omega) \in O] = \cup_{k=1}^n [\tau = k] \cap [X_k \in O]$$

To have this in \mathcal{F}_n , one needs $[\tau = k] \in \mathcal{F}_k$ for each $k \leq n$. That is $[\tau \leq k] \in \mathcal{F}_k$. Now once this is done, $[\tau > n] = [\tau \leq n]^C \in \mathcal{F}_n$ also. This motivates the following definition and shows that the requirement $[\tau \leq n] \in \mathcal{F}_n$ implies that $\omega \rightarrow X_{n \wedge \tau(\omega)}(\omega)$ is \mathcal{F}_n measurable when X_n is \mathcal{F}_n measurable and is exactly what is needed for this to happen.

Definition 29.4.1 Let (Ω, \mathcal{F}, P) be a probability space and let $\{\mathcal{F}_n\}_{n=1}^{\infty}$ be an increasing sequence of σ algebras each contained in \mathcal{F} . A stopping time is a measurable function τ which maps Ω to \mathbb{N} ,

$$\tau^{-1}(A) \in \mathcal{F} \text{ for all } A \in \mathcal{P}(\mathbb{N}),$$

such that for all $n \in \mathbb{N}$,

$$[\tau \leq n] \in \mathcal{F}_n.$$

Note this is equivalent to saying $[\tau = n] \in \mathcal{F}_n$ because $[\tau = n] = [\tau \leq n] \setminus [\tau \leq n-1]$. For τ a stopping time define \mathcal{F}_τ as follows.

$$\mathcal{F}_\tau \equiv \{A \in \mathcal{F} : A \cap [\tau \leq n] \in \mathcal{F}_n \text{ for all } n \in \mathbb{N}\} \quad (29.10)$$

These sets in \mathcal{F}_τ are referred to as “prior” to τ .

Lemma 29.4.2 The requirement 29.10 is equivalent to saying that $A \cap [\tau = n] \in \mathcal{F}_n$ for all $n \in \mathbb{N}$.

Proof: $[\tau = n] \cap A = [\tau \leq n] \cap A \setminus [\tau \leq n-1] \cap A \in \mathcal{F}_n$ so if 29.10 holds, then $A \cap [\tau = n] \in \mathcal{F}_n$ for all n . Conversely, if $A \cap [\tau = n] \in \mathcal{F}_n$, then $A \cap [\tau \leq n] = \cup_{k \leq n} A \cap [\tau = k]$ which is the union of sets in \mathcal{F}_n since the \mathcal{F}_k are increasing in k . ■

Another good thing to observe is the following lemma.

Lemma 29.4.3 If $\sigma \leq \tau$ and σ, τ are two stopping times, then $\mathcal{F}_\sigma \subseteq \mathcal{F}_\tau$.

Proof: Say $A \in \mathcal{F}_\sigma$ which means that $A \cap [\sigma \leq n] \in \mathcal{F}_n$ for all n . Now consider $A \cap [\tau = n]$. Is this in \mathcal{F}_n ? Since $\sigma \leq \tau$,

$$A \cap [\tau = n] = [\tau = n] \cap \bigcup_{j=1}^n A \cap [\sigma = j] \in \mathcal{F}_n$$

since each $[\sigma = j] \cap A \in \mathcal{F}_n$. ■

Next is a significant observation that a stochastic process $A(n)$ where $A(n)$ is \mathcal{F}_n measurable satisfies $A(\tau) \in \mathcal{F}_\tau$.

Proposition 29.4.4 Let $A(k)$ be \mathcal{F}_k measurable for each k where \mathcal{F}_k is an increasing sequence of σ algebras. Then $A(\tau)$ is \mathcal{F}_τ measurable if τ is a stopping time corresponding to the \mathcal{F}_k . If σ, τ are two stopping times, then so are $\tau \wedge \sigma$ and $\tau \vee \sigma$, the minimum and maximum of the two stopping times.

Proof: I need to show that for O an open set $[A(\tau) \in O]$ is \mathcal{F}_τ measurable. I need to show that $[A(\tau) \in O] \cap [\tau \leq k] \in \mathcal{F}_k$ for each k . It suffices to show that for each j , $[\tau = j] \cap [A(j) \in O] \cap [\tau \leq k] \in \mathcal{F}_k$ for each k . If $j \leq k$, then left side reduces to $[\tau = j] \cap [A(j) \in O]$. However, $A(j)$ is \mathcal{F}_j measurable and so $[A(j) \in O] \in \mathcal{F}_j \subseteq \mathcal{F}_k$ while $[\tau = j] \in \mathcal{F}_j \subseteq \mathcal{F}_k$ so all is well if $j \leq k$. However, if $j > k$, then the expression $[\tau = j] \cap [A(j) \leq r] \cap [\tau \leq k] = \emptyset \in \mathcal{F}_k$ and so it works in this case also. Thus $A(\tau)$ is indeed \mathcal{F}_τ measurable.

For the last claim, $[\tau \wedge \sigma \leq j] = [\tau \leq j] \cup [\sigma \leq j] \in \mathcal{F}_j$ and $[\tau \vee \sigma \leq j] = [\tau \leq j] \cap [\sigma \leq j] \in \mathcal{F}_j$. ■

Example 29.4.5 As an example of a stopping time, let X_n be \mathcal{F}_n measurable where the \mathcal{F}_n are increasing σ algebras. Let O be a Borel set, and let $\tau(\omega)$ be the first n such that $X_n(\omega) \in O$. If $X_n^{-1}(O) = \emptyset$ for all n , then $\tau(\omega) \equiv \infty$ and we consider \mathcal{F}_∞ to be \mathcal{F} . This is an example of a stopping time called the first hitting time. With these discreet processes, it is enough to let O be Borel.

Lemma 29.4.6 The first hitting time of a Borel set O is a stopping time. Also \mathcal{F}_τ is a σ algebra.

Proof: $[\tau \leq n] = \cup_{k=1}^n [\tau = k] = \cup_{k=1}^n [X_k \in O] \cap (\cup_{j < k} [X_j \in O^C])$. Now

$$[X_k \in O] \cap \left(\cup_{j < k} [X_j \in O^C] \right) \in \mathcal{F}_n$$

and so this is indeed a stopping time, being the union of finitely many sets of \mathcal{F}_n . As just noted, this means that $X_{n \wedge \tau}$ is \mathcal{F}_n measurable.

For the claim about \mathcal{F}_τ , it is obvious that Ω, \emptyset are in \mathcal{F}_τ . Suppose $A \in \mathcal{F}_\tau$. Then for k arbitrary, $A^C \cap [\tau \leq k] \cup A \cap [\tau \leq k] = [\tau \leq k]$ and so $A^C \cap [\tau \leq k] \in \mathcal{F}_k$. It is even more obvious that \mathcal{F}_τ is closed with respect to countable unions. ■

Of course τ has values i , in a countable well ordered set of numbers, $i \leq i+1$. We have the following about the relation with stopping times and conditional expectations.

Lemma 29.4.7 *Let X be in $L^1(\Omega)$. Then*

1. $\mathcal{F}_\tau \cap [\tau = i] = \mathcal{F}_i \cap [\tau = i]$ and $E(X|\mathcal{F}_\tau) = E(X|\mathcal{F}_i)$ a.e. on the set $[\tau = i]$. Also if $A \in \mathcal{F}_\tau$ or \mathcal{F}_i , then $A \cap [\tau = i] \in \mathcal{F}_i \cap \mathcal{F}_\tau$.
2. $E(X|\mathcal{F}_\tau) = E(X|\mathcal{F}_i)$ a.e. on the set $[\tau \leq i]$.

Proof: 1.) A typical set in $\mathcal{F}_\tau \cap [\tau = i]$ is $B \equiv A \cap [\tau = i]$ where $A \in \mathcal{F}_\tau$. Thus $A \cap [\tau = i] = B \in \mathcal{F}_i$ so $A \cap [\tau = i] = A \cap [\tau = i] \cap [\tau = i] = B \cap [\tau = i] \in \mathcal{F}_i \cap [\tau = i]$.

A typical set in $\mathcal{F}_i \cap [\tau = i]$ is $A \cap [\tau = i]$ where $A \in \mathcal{F}_i$. Then $A \cap [\tau = i] \cap [\tau = j] \in \mathcal{F}_j$ for all j . If $j \neq i$, you get \emptyset and if $j = i$, you get $A \cap [\tau = i] \in \mathcal{F}_i = \mathcal{F}_j$ so $A \cap [\tau = i] = B \in \mathcal{F}_\tau$ and so $A \cap [\tau = i] \cap [\tau = i] = B \cap [\tau = i] \in \mathcal{F}_\tau \cap [\tau = i]$.

For $A \in \mathcal{F}_\tau$, $A \cap [\tau = i] \in \mathcal{F}_i$ by definition. However, it is also the case, from what was just shown that $A \cap [\tau = i] \in \mathcal{F}_\tau$ because $A \cap [\tau = i] \cap [\tau = j] \in \mathcal{F}_j$ for every j . Also, if $A \in \mathcal{F}_i$, then $A \cap [\tau = i] \in \mathcal{F}_i$ by definition of a stopping time and $A \cap [\tau = i] \cap [\tau = j] \in \mathcal{F}_j$ for every j . Thus if A is either in \mathcal{F}_τ or \mathcal{F}_i , then $[\tau = i] \cap A \in \mathcal{F}_i \cap \mathcal{F}_\tau$.

Now let $A \in \mathcal{F}_\tau$. Then

$$\int_{A \cap [\tau = i]} E(X|\mathcal{F}_i) dP = \int_{A \cap [\tau = i]} X dP \equiv \int_{A \cap [\tau = i]} E(X|\mathcal{F}_\tau) dP$$

2.) If $A \in \mathcal{F}_\tau$, then $A \cap [\tau = j] \in \mathcal{F}_j$ by definition of \mathcal{F}_τ and so by definition of conditional expectation,

$$\begin{aligned} \int_{A \cap [\tau \leq i]} E(X|\mathcal{F}_i) dP &= \sum_{j=1}^i \int_{A \cap [\tau = j]} E(X|\mathcal{F}_i) dP = \sum_{j=1}^i \int_{A \cap [\tau = j]} E(X|\mathcal{F}_j) dP \\ &= \sum_{j=1}^i \int_{A \cap [\tau = j]} E(X|\mathcal{F}_\tau) dP = \int_{A \cap [\tau \leq i]} E(X|\mathcal{F}_\tau) dP \end{aligned}$$

Since $A \in \mathcal{F}_\tau$ is arbitrary, it follows that $E(X|\mathcal{F}_i) = E(X|\mathcal{F}_\tau)$ a.e. on the set $[\tau \leq i]$. ■

29.4.1 Optional Sampling for Martingales

Now it is time for the optional sampling theorem. Suppose $\{X_n\}$ is a martingale. In particular, each $X_n \in L^1(\Omega; W)$ and $E(X_n | \mathcal{F}_k) = X_k$ whenever $k \leq n$. We can assume X_n has values in some separable Banach space. Then $\|X_n\|$ is a sub-martingale because if $k \leq n$, then if $A \in \mathcal{F}_k$,

$$\int_A E(\|X_n\| | \mathcal{F}_k) dP \geq \int_A \|E(X_n | \mathcal{F}_k)\| dP = \int_A \|X_k\| dP$$

Now suppose we have two stopping times τ and σ and τ is bounded meaning it has values in $\{1, 2, \dots, n\}$.

The optional sampling theorem says the following. For $M(n)$ a martingale, $\|M(n)\| \in L^1$, then the following holds a.e.

$$M(\sigma \wedge \tau) = E(M(\tau) | \mathcal{F}_\sigma)$$

Furthermore, it all makes sense. First of all, why does it make sense? We need to verify that $M(\tau)$ is integrable.

$$\int \|M(\tau)\| = \sum_{k=1}^n \int_{[\tau=k]} \|M(k)\| dP < \infty$$

Similarly, $M(\sigma \wedge \tau)$ is integrable. The reason is that $\sigma \wedge \tau$ is a stopping time which is bounded by n . Thus the above follows with τ replaced with $\sigma \wedge \tau$. Why is $\sigma \wedge \tau$ a stopping time? It is because $[\sigma \wedge \tau \leq k] = [\sigma \leq k] \cup [\tau \leq k] \in \mathcal{F}_k$. It is also clear that $\tau = i \in \mathbb{N}$ will be a stopping time.

Now let $A \in \mathcal{F}_\sigma$. Then using Lemma 29.4.7 as needed,

$$\int_A M(\sigma \wedge \tau) = \sum_{i=1}^n \int_{A \cap [\tau=i]} M(\sigma \wedge i) = \sum_{i=1}^n \sum_{j=1}^{\infty} \int_{A \cap [\tau=i] \cap [\sigma=j]} M(j \wedge i)$$

If $j \leq i$,

$$M(j \wedge i) = M(j) = E(M(i) | \mathcal{F}_j).$$

If $j > i$,

$$M(j \wedge i) = M(i) = E(M(i) | \mathcal{F}_j).$$

On $[j = \sigma]$, $E(M(i) | \mathcal{F}_j) = E(M(i) | \mathcal{F}_\sigma)$ Thus the last term in the above is

$$= \sum_{i=1}^n \sum_{j=1}^{\infty} \int_{A \cap [\tau=i] \cap [\sigma=j]} E(M(i) | \mathcal{F}_\sigma) = \sum_{i=1}^n \int_{A \cap [\tau=i]} E(M(i) | \mathcal{F}_\sigma)$$

Now $\mathcal{X}_{A \cap [\tau=i]} M(i) = \mathcal{X}_{A \cap [\tau=i]} M(\tau)$ so

$$\begin{aligned} &= \sum_{i=1}^n \int E(\mathcal{X}_{A \cap [\tau=i]} M(i) | \mathcal{F}_\sigma) = \sum_{i=1}^n \int E(\mathcal{X}_{A \cap [\tau=i]} M(\tau) | \mathcal{F}_\sigma) \\ &= \int E(\mathcal{X}_A M(\tau) | \mathcal{F}_\sigma) = \int_A E(M(\tau) | \mathcal{F}_\sigma) \end{aligned}$$

Since A is an arbitrary element of \mathcal{F}_σ , this shows the optional sampling theorem that $M(\sigma \wedge \tau) = E(M(\tau) | \mathcal{F}_\sigma)$.

Proposition 29.4.8 *Let M be a martingale having values in some separable Banach space. Let τ be a bounded stopping time and let σ be another stopping time. Then everything makes sense in the following formula and*

$$M(\sigma \wedge \tau) \equiv M_{\sigma \wedge \tau} = E(M(\tau) | \mathcal{F}_\sigma) \text{ a.e.}$$

29.4.2 Optional Sampling Theorem for Sub-Martingales

What about the case where $\{X(n)\}$ is a sub-martingale? Shouldn't there be something like the conclusion of Proposition 29.4.8? This requires a very interesting theorem which involves the decomposition of a sub-martingale into a sum. Recall $\{X(k)\}_{k=1}^\infty$ is a sub-martingale if

$$E(X(k+1) | \mathcal{F}_k) \geq X(k)$$

where the \mathcal{F}_k are an increasing sequence of σ algebras in the usual way. The following is the very interesting result about writing a sub-martingale as the sum of an increasing process and a martingale.

Lemma 29.4.9 *Let $\{X(k)\}_{k=0}^\infty$ be a sub-martingale adapted to the increasing sequence of σ algebras, $\{\mathcal{F}_k\}$. Then there exists a unique increasing process $\{A(k)\}_{k=0}^\infty$ such that $A(0) = 0$ and $A(k+1)$ is \mathcal{F}_k measurable for all k and a martingale, $\{M(k)\}_{k=0}^\infty$ such that*

$$X(k) = A(k) + M(k).$$

Furthermore, for τ a stopping time, $A(\tau)$ is \mathcal{F}_τ measurable.

Proof: Define $\sum_{k=0}^{-1} \equiv 0$. First consider the uniqueness assertion. Suppose A is a process which does what it is supposed to do.

$$\begin{aligned} \sum_{k=0}^{n-1} E(X(k+1) - X(k) | \mathcal{F}_k) &= \sum_{k=0}^{n-1} E(A(k+1) - A(k) | \mathcal{F}_k) \\ &+ \sum_{k=0}^{n-1} E(M(k+1) - M(k) | \mathcal{F}_k) \end{aligned}$$

Then since $\{M(k)\}$ is a martingale,

$$\sum_{k=0}^{n-1} E(X(k+1) - X(k) | \mathcal{F}_k) = \sum_{k=0}^{n-1} A(k+1) - A(k) = A(n)$$

This shows uniqueness and gives a formula for $A(n)$ assuming it exists. It is only a matter of verifying this does work. Define

$$A(n) \equiv \sum_{k=0}^{n-1} E(X(k+1) - X(k) | \mathcal{F}_k), \quad A(0) = 0.$$

Then A is increasing because from the definition,

$$A(n+1) - A(n) = E(X(n+1) - X(n) | \mathcal{F}_n) \geq 0.$$

Also from the definition above, $A(n)$ is \mathcal{F}_{n-1} measurable, so consider

$$\{X(k) - A(k)\}.$$

Why is this a martingale?

$$\begin{aligned}
 E(X(k+1) - A(k+1) | \mathcal{F}_k) &= E(X(k+1) | \mathcal{F}_k) - A(k+1) \\
 &= E(X(k+1) | \mathcal{F}_k) - \sum_{j=0}^k E(X(j+1) - X(j) | \mathcal{F}_j) \\
 &= E(X(k+1) | \mathcal{F}_k) - E(X(k+1) - X(k) | \mathcal{F}_k) \\
 &\quad - \sum_{j=0}^{k-1} E(X(j+1) - X(j) | \mathcal{F}_j) \\
 &= X(k) - \sum_{j=0}^{k-1} E(X(j+1) - X(j) | \mathcal{F}_j) = X(k) - A(k)
 \end{aligned}$$

Let $M(k) \equiv X(k) - A(k)$. $A(\tau)$ is \mathcal{F}_τ measurable by Proposition 29.4.4. ■

Note the nonnegative integers could be replaced with any finite set or ordered countable set of numbers with no change in the conclusions of this lemma or the above optional sampling theorem.

Next consider the case of a sub-martingale.

Theorem 29.4.10 *Let $\{X(k)\}$ be a sub-martingale with respect to the increasing sequence of σ algebras, $\{\mathcal{F}_k\}$ and let σ, τ be two stopping times such that τ is bounded. Then $X(\tau)$ defined as*

$$\omega \rightarrow X(\tau(\omega))$$

is integrable and

$$X(\sigma \wedge \tau) \leq E(X(\tau) | \mathcal{F}_\sigma).$$

Proof: The claim about $X(\tau)$ being integrable is the same as done earlier. If $\tau \leq l$,

$$E(|X(\tau(\omega))|) = \sum_{i=1}^l \int_{[\tau=i]} |X(i)| dP < \infty$$

By Lemma 29.4.9 there is a martingale, $\{M(k)\}$ and an increasing process $\{A(k)\}$ such that $A(k+1)$ is \mathcal{F}_k measurable such that

$$X(k) = M(k) + A(k).$$

Then from the fact A is increasing,

$$\begin{aligned}
 E(X(\tau) | \mathcal{F}_\sigma) &= E(M(\tau) + A(\tau) | \mathcal{F}_\sigma) = M(\tau \wedge \sigma) + E(A(\tau) | \mathcal{F}_\sigma) \\
 &\geq M(\tau \wedge \sigma) + E(A(\tau \wedge \sigma) | \mathcal{F}_\sigma) \\
 &= M(\tau \wedge \sigma) + A(\tau \wedge \sigma) \equiv X(\tau \wedge \sigma).
 \end{aligned}$$

because in the above, it follows from Lemma 29.4.9, $A(\tau \wedge \sigma)$ is $\mathcal{F}_{\tau \wedge \sigma}$ measurable and from Lemma 29.4.3, $\mathcal{F}_{\tau \wedge \sigma} \subseteq \mathcal{F}_\sigma$ and so $E(A(\tau \wedge \sigma) | \mathcal{F}_\sigma) = A(\tau \wedge \sigma)$. ■

Say τ is bounded by n and σ is a stopping time. A useful way to remember the above theorem is in the following proposition.

Proposition 29.4.11 *Suppose τ is a bounded stopping time and σ is a stopping time. Then if $\{X(k)\}$ is a sub-martingale, then $X(1), X(\tau \wedge \sigma), X(\tau)$ is also a sub-martingale.*

Proof: Consider $X(1), X(\tau \wedge \sigma), X(\tau)$. Then

$$E(X(\tau) | \mathcal{F}_{\tau \wedge \sigma}) \geq X(\sigma \wedge \tau \wedge \tau) = X(\sigma \wedge \tau).$$

Also $E(X(\sigma \wedge \tau) | \mathcal{F}_1) \geq X(\sigma \wedge \tau \wedge 1) = X(1)$ so this stochastic process is a sub-martingale. ■

This optional sampling theorem gives a convenient way to consider the Doob maximal estimate presented earlier.

Proposition 29.4.12 *Let $\{X(k)\}$ be a real valued sub-martingale, and let $\lambda > 0$. Then for $X_n^* \equiv \max\{X_k : k \leq n\}$ as earlier,*

$$\int_{[X_n^* \leq \lambda]} X(n)^+ dP \geq \lambda P([X_n^* > \lambda]) = \lambda P([(X_n^*)^+ > \lambda])$$

Proof: Let $\tau = n$ and let σ be the first hitting time of the set (λ, ∞) by $X(k)$. Then $\omega \in [X_n^* > \lambda]$ if and only if for some $k \leq n, X_k > \lambda$ if and only if $\sigma(\omega) = k$ for some $k \leq n$. By Proposition 29.4.11 $X(1), X(\sigma \wedge n), X(n)$ is a sub-martingale and the set of interest $[X_n^* > \lambda]$ is the one where $\sigma < \infty$. Then

$$\begin{aligned} E(X(n)) &\geq E(X(\sigma \wedge n)) = \int_{[\sigma < \infty]} X(\sigma \wedge n) dP + \int_{[\sigma = \infty]} X(\sigma \wedge n) dP \\ &= \int_{[\sigma < \infty]} X(\sigma \wedge n) dP + \int_{[\sigma = \infty]} X(n) dP \\ &\geq \lambda P([\sigma < \infty]) + \int_{[\sigma = \infty]} X(n) dP \end{aligned}$$

Therefore,

$$E(\mathcal{X}_{[X_n^* > \lambda]} X(n)^+) \geq E(\mathcal{X}_{[X_n^* > \lambda]} X(n)) \geq \lambda P([\sigma < \infty]) = \lambda P([X_n^* > \lambda]). \blacksquare$$

Now let $X_{n*} = \min\{X(k) : k \leq n\}$. What about $P([X_{n*} < -\lambda])$? Let σ be the first hitting time for $(-\infty, -\lambda)$ and note that $[X_{n*} < -\lambda]$ consists of the set of ω where $\sigma(\omega) < \infty$. As noted in Proposition 29.4.11, $X(1), X(\sigma \wedge n), X(n)$ is a sub-martingale. Thus

$$\begin{aligned} \int X(1) dP &\leq \int_{[\sigma < \infty]} X(\sigma \wedge n) dP + \int_{[\sigma = \infty]} X(\sigma \wedge n) dP \\ &= \int_{[\sigma < \infty]} X(\sigma \wedge n) dP + \int_{[\sigma = \infty]} X(n) dP \end{aligned}$$

It follows that

$$\int X(1) dP - \int_{[\sigma = \infty]} X(n) dP \leq \int_{[\sigma < \infty]} X(\sigma \wedge n) dP \leq -\lambda P([\sigma < \infty])$$

and so,

$$\lambda P([\sigma < \infty]) = \lambda P([X_{n*} < -\lambda]) \leq \int |X(1)| + |X(n)| dP$$

Therefore, we obtain the following theorem which is a maximal estimate for sub-martingales.

Theorem 29.4.13 *Let $\{X(k)\}$ be a real sub-martingale and let $\lambda > 0$ be given. Then*

$$P([\max\{|X_k|, k = 1, \dots, n\} > \lambda]) \leq \frac{2}{\lambda} \int |X(1)| + |X(n)| dP$$

Proof: $[\max\{|X_k|, k = 1, \dots, n\} > \lambda] \subseteq [X_n^* > \lambda] \cup [X_{n*} < -\lambda]$ and so

$$\begin{aligned} P([\max\{|X_k|, k = 1, \dots, n\} > \lambda]) &\leq \frac{1}{\lambda} \int X_n^+ dP + \frac{1}{\lambda} \int (|X(1)| + |X(n)|) dP \\ &\leq \frac{2}{\lambda} \int |X(1)| + |X(n)| dP. \blacksquare \end{aligned}$$

29.5 Reverse Sub-martingale Convergence Theorem

Sub-martingale: $E(X_{n+1} | \mathcal{F}_n) \geq X_n$. Reverse sub-martingale: $E(X_n | \mathcal{F}_{n+1}) \geq X_{n+1}$ and here the \mathcal{F}_n are decreasing.

Definition 29.5.1 *Let $\{X_n\}_{n=0}^\infty$ be a sequence of real random variables such that $E(|X_n|) < \infty$ for all n and let $\{\mathcal{F}_n\}$ be a sequence of σ algebras such that $\mathcal{F}_n \supseteq \mathcal{F}_{n+1}$ for all n . Then $\{X_n\}$ is called a reverse sub-martingale if for all n ,*

$$E(X_n | \mathcal{F}_{n+1}) \geq X_{n+1}.$$

Note it is just like a sub-martingale only the indices and σ algebras are going the other way. Here is an interesting lemma. This lemma gives uniform integrability for a reverse sub-martingale. The application I have in mind in the next lemma is that $\sup_n E(|X_n|) < \infty$ but it is stated more generally and this condition appears to be obtained for free given the existence of X_∞ in the following lemma.

Lemma 29.5.2 *Suppose for each n , $E(|X_n|) < \infty$, X_n is \mathcal{F}_n measurable, $\mathcal{F}_{n+1} \subseteq \mathcal{F}_n$ for all $n \in \mathbb{N}$, and there exist X_∞ \mathcal{F}_∞ measurable such that $\mathcal{F}_\infty \subseteq \mathcal{F}_n$ for all n and X_0 \mathcal{F}_0 measurable such that $\mathcal{F}_0 \supseteq \mathcal{F}_n$ for all n such that for all $n \in \{0, 1, \dots\}$,*

$$E(X_n | \mathcal{F}_{n+1}) \geq X_{n+1}, \quad E(X_n | \mathcal{F}_\infty) \geq X_\infty,$$

where $E(|X_\infty|) < \infty$. Then $\{X_n : n \in \mathbb{N}\}$ is equi-integrable.

Proof:

$$E(X_{n+1}) \leq E(E(X_n | \mathcal{F}_{n+1})) = E(X_n)$$

Therefore, the sequence $\{E(X_n)\}$ is a decreasing sequence bounded below by $E(X_\infty)$ so it has a limit. I am going to show the functions are equi-integrable. Let k be large enough that

$$\left| E(X_k) - \lim_{m \rightarrow \infty} E(X_m) \right| < \varepsilon \quad (29.11)$$

and suppose $n > k$. Then if $\lambda > 0$,

$$\int_{[|X_n| \geq \lambda]} |X_n| dP = \int_{[X_n \geq \lambda]} X_n dP + \int_{[X_n \leq -\lambda]} (-X_n) dP$$

$$\begin{aligned}
&= \int_{[X_n \geq \lambda]} X_n dP + \int_{\Omega} (-X_n) dP + \int_{[-X_n < \lambda]} X_n dP \\
&\leq \int_{[X_n \geq \lambda]} E(X_k | \mathcal{F}_n) dP + \int_{\Omega} (-X_n) dP + \int_{[-X_n < \lambda]} E(X_k | \mathcal{F}_n) dP \\
&\leq \int_{[X_n \geq \lambda]} X_k dP + \int_{\Omega} (-X_k) dP + \int_{[X_n > -\lambda]} X_k dP + \varepsilon \\
&= \int_{[X_n \geq \lambda]} X_k dP - \left(\int_{\Omega} X_k - \int_{[X_n > -\lambda]} X_k dP \right) + \varepsilon \\
&= \int_{[X_n \geq \lambda]} X_k dP - \left(\int_{[X_n \leq -\lambda]} X_k dP \right) + \varepsilon \\
&= \int_{[X_n \geq \lambda]} X_k dP + \left(\int_{[-X_n \geq \lambda]} (-X_k) dP \right) + \varepsilon = \int_{[|X_n| \geq \lambda]} |X_k| dP + \varepsilon
\end{aligned}$$

Applying the maximal inequality for sub-martingales, Theorem 29.4.13,

$$P\left(\left[\max\{|X_j| : j = n, \dots, 1\} \geq \lambda\right]\right) \leq \frac{1}{\lambda} (E(|X_0|) + E(|X_{\infty}|)) \leq \frac{C}{\lambda}$$

and taking sup for all n , $P\left(\left[\sup\{|X_j|\} \geq \lambda\right]\right) \leq \frac{C}{\lambda}$. From the above, for $n > k$,

$$\int_{[|X_n| \geq \lambda]} |X_n| dP \leq \int_{[|X_n| \geq \lambda]} |X_k| dP + \varepsilon, \quad P(|X_n| \geq \lambda) \leq \frac{C}{\lambda}$$

Since the single function X_k is equi-integrable, it follows that for all λ large enough, $\int_{[|X_n| \geq \lambda]} |X_n| dP \leq 2\varepsilon$ for all $n > k$. Since there are only finitely many X_j for $j \leq k$, this shows $\{X_n\}$ is equi-integrable. Hence $\{X_n\}$ is uniformly integrable. ■

Note that this also gives X_n bounded in $L^1(\Omega)$ from Proposition 10.9.6 on Page 293. Now with this lemma and the upcrossing lemma it is easy to prove an important convergence theorem.

Theorem 29.5.3 *Let $\{X_n, \mathcal{F}_n\}_{n=0}^{\infty}$ be a backwards sub-martingale as described above and suppose $\sup_{n \geq 0} E(|X_n|) < \infty$. Then $\{X_n\}$ converges a.e. and in $L^1(\Omega)$ to a function, $X_{\infty} \in L^1(\Omega)$.*

Proof: By the upcrossing lemma applied to the sub-martingale $\{X_k\}_{k=0}^N$, the number of upcrossings (Downcrossings is probably a better term. They are upcrossings as n gets smaller.) of the interval $[a, b]$ satisfies the inequality

$$E\left(U_{[a,b]}^N\right) \leq \frac{1}{b-a} C$$

Letting $N \rightarrow \infty$, it follows the expected number of upcrossings, $E(U_{[a,b]})$ is bounded. Therefore, there exists a set of measure 0 N_{ab} such that if $\omega \notin N_{ab}$, $U_{[a,b]}(\omega) < \infty$. Let $N = \cup\{N_{ab} : a, b \in \mathbb{Q}\}$. Then for $\omega \notin N$,

$$\limsup_{n \rightarrow \infty} X_n(\omega) = \liminf_{n \rightarrow \infty} X_n(\omega)$$

because if inequality holds, then letting

$$\liminf_{n \rightarrow \infty} X_n(\omega) < a < b < \limsup_{n \rightarrow \infty} X_n(\omega)$$

it would follow $U_{[a,b]}(\omega) = \infty$, contrary to $\omega \notin N_{ab}$.

Let $X_\infty(\omega) \equiv \lim_{n \rightarrow \infty} X_n(\omega)$. Then by Fatou's lemma,

$$\int_{\Omega} |X_\infty(\omega)| dP \leq \liminf_{n \rightarrow \infty} \int_{\Omega} |X_n| dP < \infty.$$

and so X_∞ is in $L^1(\Omega)$. By the Vitali convergence theorem and Lemma 29.5.2 which shows $\{|X_n|\}$ is uniformly integrable, it follows

$$\lim_{n \rightarrow \infty} \int_{\Omega} |X_\infty(\omega) - X_n(\omega)| dP = 0. \blacksquare$$

29.6 Strong Law of Large Numbers

There is a version of the strong law of large numbers which does not depend on the random variables having finite variance. First are some preparatory lemmas. The approach followed here is from Ash [3].

The message of the following lemma, $E(X_k | \sigma(S_n)) = E(X_k | \sigma(S_n, \mathbf{Y}))$ makes sense. The expectation of X_k for $k \leq n$ given the value of $S_n \equiv \sum_{k=1}^n X_k$ is the same as the expectation of X_k given the value of S_n and X_{n+1}, \dots . It makes intuitive sense because the random variables are independent so knowledge of X_j for $j \geq n+1$ should be irrelevant to the expectation of X_k .

Lemma 29.6.1 *Let $\{X_n\}$ be a sequence of independent random variables such that $E(|X_k|) < \infty$ for all k and let $S_n \equiv \sum_{k=1}^n X_k$. Then for $k \leq n$,*

$$E(X_k | \sigma(S_n)) = E(X_k | \sigma(S_n, \mathbf{Y})) \text{ a.e.} \quad (29.12)$$

where $\mathbf{Y} = (X_{n+1}, X_{n+2}, \dots) \in \mathbb{R}^{\mathbb{N}} \equiv \prod_{i=1}^{\infty} \mathbb{R}$. Also for $k \leq n$ as above,

$$\sigma(S_n, \mathbf{Y}) = \sigma(S_n, S_{n+1}, \dots).$$

Proof: Note that both X_k and S_n are measurable with respect to $\sigma(X_1, \dots, X_n)$ and $\sigma(X_1, \dots, X_n)$ and $\sigma(\mathbf{Y})$ are independent. Therefore, by Proposition 29.2.1 29.12 holds.

It only remains to prove the last assertion. For $k > 0$,

$$X_{n+k} = S_{n+k} - S_{n+k-1}$$

Thus

$$\begin{aligned} \sigma(S_n, \mathbf{Y}) &= \sigma(S_n, X_{n+1}, \dots) \\ &= \sigma(S_n, (S_{n+1} - S_n), (S_{n+2} - S_{n+1}), \dots) \end{aligned}$$

Thus, by induction, each S_{n+k} is measurable with respect to $\sigma(S_n, \mathbf{Y})$ and so,

$$\sigma(S_n, S_{n+1}, \dots) \subseteq \sigma(S_n, \mathbf{Y})$$

To get the other inclusion,

$$\sigma(S_n, S_{n+1}, \dots) = \sigma(S_n, X_{n+1} + S_n, X_{n+2} + S_{n+1}, \dots)$$

so by induction, each X_{n+k} and S_n is measurable with respect to $\sigma(S_n, S_{n+1}, \dots)$ and so $\sigma(S_n, \mathbf{Y}) \subseteq \sigma(S_n, S_{n+1}, \dots)$. \blacksquare

Lemma 29.6.2 *Let $\{X_k\}$ be a sequence of independent identically distributed random variables such that $E(|X_k|) < \infty$. Then letting $S_n = \sum_{k=1}^n X_k$, it follows that for $k \leq n$*

$$E(X_k | \sigma(S_n, S_{n+1}, \dots)) = E(X_k | \sigma(S_n)) = \frac{S_n}{n}.$$

Proof: It was shown in Lemma 29.6.1 the first equality holds. It remains to show the second. Letting $A = S_n^{-1}(B)$ where B is Borel, it follows there exists $B' \subseteq \mathbb{R}^n$ a Borel set such that

$$S_n^{-1}(B) = (X_1, \dots, X_n)^{-1}(B').$$

Then

$$\begin{aligned} \int_A E(X_k | \sigma(S_n)) dP &= \int_{S_n^{-1}(B)} X_k dP \\ &= \int_{(X_1, \dots, X_n)^{-1}(B')} X_k dP = \int_{(X_1, \dots, X_n)^{-1}(B')} x_k d\lambda_{(X_1, \dots, X_n)} \\ &= \int \cdots \int \mathcal{X}_{(X_1, \dots, X_n)^{-1}(B')}(x) x_k d\lambda_{X_1} d\lambda_{X_2} \cdots d\lambda_{X_n} \\ &= \int \cdots \int \mathcal{X}_{(X_1, \dots, X_n)^{-1}(B')}(x) x_l d\lambda_{X_1} d\lambda_{X_2} \cdots d\lambda_{X_n} \\ &= \int_A E(X_l | \sigma(S_n)) dP \end{aligned}$$

and so since $A \in \sigma(S_n)$ is arbitrary,

$$E(X_l | \sigma(S_n)) = E(X_k | \sigma(S_n))$$

for each $k, l \leq n$. Therefore,

$$S_n = E(S_n | \sigma(S_n)) = \sum_{j=1}^n E(X_j | \sigma(S_n)) = nE(X_k | \sigma(S_n)) \text{ a.e.}$$

and so

$$E(X_k | \sigma(S_n)) = \frac{S_n}{n} \text{ a.e.}$$

as claimed. ■

With this preparation, here is the strong law of large numbers for identically distributed random variables.

Theorem 29.6.3 *Let $\{X_k\}$ be a sequence of independent identically distributed random variables such that $E(|X_k|) < \infty$ for all k . Since these are identically distributed, $E(|X_k|)$ does not depend on k and so the process is bounded. Letting $m = E(X_k)$,*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n X_k(\omega) = m \text{ a.e.}$$

and convergence also takes place in $L^1(\Omega)$.

Proof: Consider the reverse sub-martingale $\{E(X_1 | \sigma(S_n, S_{n+1}, \dots))\}$. By Theorem 29.5.3, this converges a.e. and in $L^1(\Omega)$ to a random variable X_∞ . However, from Lemma 29.6.2, $E(X_1 | \sigma(S_n, S_{n+1}, \dots)) = S_n/n$. Therefore, S_n/n converges a.e. and in $L^1(\Omega)$ to X_∞ . I need to argue that X_∞ is constant and also that it equals m . For $a \in \mathbb{R}$ let

$$E_a \equiv [X_\infty \geq a]$$

For a small enough, $P(E_a) \neq 0$. Then since E_a is a tail event for the independent random variables, $\{X_k\}$ it follows from the Kolmogorov zero one law, Theorem 26.7.4, that $P(E_a) = 1$. Let $b \equiv \sup\{a : P(E_a) = 1\}$. The sets, E_a are decreasing as a increases. Let $\{a_n\}$ be a strictly increasing sequence converging to b . Then

$$[X_\infty \geq b] = \cap_n [X_\infty \geq a_n]$$

and so $1 = P(E_b) = \lim_{n \rightarrow \infty} P(E_{a_n})$. On the other hand, if $c > b$, then $P(E_c) < 1$ and so $P(E_c) = 0$. Hence $P([X = b]) = 1$. It remains to show $b = m$. This is easy because by the L^1 convergence,

$$b = \int_{\Omega} X_\infty dP = \lim_{n \rightarrow \infty} \int_{\Omega} \frac{S_n}{n} dP = \lim_{n \rightarrow \infty} m = m. \blacksquare$$

Chapter 30

Continuous Stochastic Processes

The change here is that the stochastic process will depend on $t \in I$ an interval rather than $n \in \mathbb{N}$. Everything becomes much more technical.

30.1 Fundamental Definitions and Properties

Here E will be a separable Banach space and $\mathcal{B}(E)$ will be the Borel sets of E . Let (Ω, \mathcal{F}, P) be a probability space and I will be an interval of \mathbb{R} . A set of E valued random variables, one for each $t \in I$, $\{X(t) : t \in I\}$ is called a stochastic process. Thus for each t , $X(t)$ is a measurable function of $\omega \in \Omega$. Set $X(t, \omega) \equiv X(t)(\omega)$. Functions $t \rightarrow X(t, \omega)$ are called trajectories. Thus there is a trajectory for each $\omega \in \Omega$. A stochastic process, Y is called a version or a modification of a stochastic process X if for all $t \in I$,

$$X(t, \omega) = Y(t, \omega) \text{ a.e. } \omega$$

There are several descriptions of stochastic processes.

1. X is measurable if $X(\cdot, \cdot) : I \times \Omega \rightarrow E$ is $B(I) \times \mathcal{F}$ measurable. Note that a stochastic process X is not necessarily measurable.
2. X is stochastically continuous at $t_0 \in I$ means: for all $\varepsilon > 0$ and $\delta > 0$ there exists $\rho > 0$ such that

$$P(\|X(t) - X(t_0)\| \geq \varepsilon) \leq \delta \text{ whenever } |t - t_0| < \rho, t \in I.$$

Note the above condition says that for each $\varepsilon > 0$,

$$\lim_{t \rightarrow t_0} P(\|X(t) - X(t_0)\| \geq \varepsilon) = 0.$$

3. X is stochastically continuous if it is stochastically continuous at every $t \in I$.
4. X is stochastically uniformly continuous if for every $\varepsilon, \delta > 0$ there exists $\rho > 0$ such that whenever $s, t \in I$ with $|s - t| < \rho$, it follows

$$P(\|X(t) - X(s)\| \geq \varepsilon) \leq \delta.$$

5. X is mean square continuous at $t_0 \in I$ if

$$\lim_{t \rightarrow t_0} E(\|X(t) - X(t_0)\|^2) \equiv \lim_{t \rightarrow t_0} \int_{\Omega} \|X(t)(\omega) - X(t_0)(\omega)\|^2 dP = 0.$$

6. X is mean square continuous in I if it is mean square continuous at every point of I .
7. X is continuous with probability 1 or continuous if $t \rightarrow X(t, \omega)$ is continuous for all ω outside some set of measure 0.
8. X is Hölder continuous if $t \rightarrow X(t, \omega)$ is Hölder continuous for a.e. ω .

Lemma 30.1.1 *A stochastically continuous process on $[a, b] \equiv I$ is uniformly stochastically continuous on $[a, b] \equiv I$.*

Proof: If this is not so, there exists $\varepsilon, \delta > 0$ and points of I, s_n, t_n such that even though $|t_n - s_n| < \frac{1}{n}$,

$$P(\|X(s_n) - X(t_n)\| \geq \varepsilon) > \delta. \quad (30.1)$$

Taking a subsequence, still denoted by s_n and t_n there exists $t \in I$ such that the above hold and $\lim_{n \rightarrow \infty} s_n = \lim_{n \rightarrow \infty} t_n = t$. Then

$$\begin{aligned} & P(\|X(s_n) - X(t_n)\| \geq \varepsilon) \\ & \leq P(\|X(s_n) - X(t)\| \geq \varepsilon/2) + P(\|X(t) - X(t_n)\| \geq \varepsilon/2). \end{aligned}$$

But the sum of the last two terms converges to 0 as $n \rightarrow \infty$ by stochastic continuity of X at t , violating 30.1 for all n large enough. ■

For a stochastically continuous process defined on a closed and bounded interval, there always exists a measurable version. This is significant because then you can do things with product measure and iterated integrals.

Proposition 30.1.2 *Let X be a stochastically continuous process defined on a closed interval, $I \equiv [a, b]$. Then there exists a measurable version of X .*

Proof: By Lemma 30.1.1 X is uniformly stochastically continuous and so there exists a sequence of positive numbers, $\{\rho_n\}$ such that if $|s - t| < \rho_n$, then

$$P\left(\left[\|X(t) - X(s)\| \geq \frac{1}{2^n}\right]\right) \leq \frac{1}{2^n}. \quad (30.2)$$

Then let $\{t_0^n, t_1^n, \dots, t_{m_n}^n\}$ be a partition of $[a, b]$ in which $|t_i^n - t_{i-1}^n| < \rho_n$. Now define X_n as follows:

$$X_n(t) \equiv \sum_{i=1}^{m_n} X(t_{i-1}^n) \mathcal{X}_{[t_{i-1}^n, t_i^n)}(t), \quad X_n(b) \equiv X(b).$$

Then X_n is obviously $B(I) \times \mathcal{F}$ measurable because it is the sum of functions which are. Consider the set A on which $\{X_n(t, \omega)\}$ is a Cauchy sequence. This set is of the form

$$A = \bigcap_{n=1}^{\infty} \bigcup_{m=1}^{\infty} \bigcap_{p, q \geq m} \left[\|X_p - X_q\| < \frac{1}{n} \right]$$

and so it is a $B(I) \times \mathcal{F}$ measurable set. Now define

$$Y(t, \omega) \equiv \begin{cases} \lim_{n \rightarrow \infty} X_n(t, \omega) & \text{if } (t, \omega) \in A \\ 0 & \text{if } (t, \omega) \notin A \end{cases}$$

I claim $Y(t, \omega) = X(t, \omega)$ for a.e. ω . To see this, consider 30.2. From the construction of X_n , it follows that for each t ,

$$P\left(\left[\|X_n(t) - X(t)\| \geq \frac{1}{2^n}\right]\right) \leq \frac{1}{2^n} \quad (30.3)$$

Also, for a fixed t , if $X_n(t, \omega)$ fails to converge to $X(t, \omega)$, then ω must be in infinitely many of the sets,

$$B_n \equiv \left[\|X_n(t) - X(t)\| \geq \frac{1}{2^n} \right]$$

which is a set of measure zero by the Borel Cantelli lemma and 30.3. Recall why this is so.

$$P(\cap_{k=1}^{\infty} \cup_{n=k}^{\infty} B_n) \leq \sum_{n=k}^{\infty} P(B_n) < \frac{1}{2^{k-1}}$$

Therefore, for each $t, (t, \omega) \in A$ for a.e. ω . Hence $X(t) = Y(t)$ a.e. and so Y is a measurable version of X . ■

One also has the following lemma about extending a process from a dense subset.

Lemma 30.1.3 *Let D be a dense subset of an interval, $I = [0, T]$ and suppose $X : D \rightarrow E$ satisfies*

$$\|X(d) - X(d')\| \leq C|d - d'|^\gamma$$

for all $d', d \in D$. Then X extends uniquely to a continuous Y defined on $[0, T]$ such that

$$\|Y(t) - Y(t')\| \leq C|t - t'|^\gamma.$$

Proof: Let $t \in I$ and let $d_k \rightarrow t$ where $d_k \in D$. Then $\{X(d_k)\}$ is a Cauchy sequence because $\|X(d_k) - X(d_m)\| \leq C|d_k - d_m|^\gamma$. Therefore, $X(d_k)$ converges. The thing it converges to will be called $Y(t)$. Note this is well defined, giving $X(t)$ if $t \in D$. Also, if $d_k \rightarrow t$ and $d'_k \rightarrow t$, then $\|X(d_k) - X(d'_k)\| \leq C|d_k - d'_k|^\gamma$ and so $X(d_k)$ and $X(d'_k)$ converge to the same thing. Therefore, it makes sense to define $Y(t) \equiv \lim_{d \rightarrow t} X(d)$. It only remains to verify the estimate. But letting $|d - t|$ and $|d' - t'|$ be small enough,

$$\begin{aligned} \|Y(t) - Y(t')\| &= \|X(d) - X(d')\| + \varepsilon \\ &\leq C|d' - d| + \varepsilon \leq C|t - t'| + 2\varepsilon. \end{aligned}$$

Since ε is arbitrary, this proves the existence part of the lemma. Uniqueness follows from observing that $Y(t)$ must equal $\lim_{d \rightarrow t} X(d)$. ■

30.2 Kolmogorov Čentsov Continuity Theorem

Lemma 30.2.1 *Let r_j^m denote $j(\frac{T}{2^m})$ where $j \in \{0, 1, \dots, 2^m\}$. Also let $D_m = \{r_j^m\}_{j=1}^{2^m}$ and $D = \cup_{m=1}^{\infty} D_m$. Suppose $X(t)$ satisfies*

$$\|X(r_{j+1}^k) - X(r_j^k)\| \leq 2^{-\gamma k} \quad (30.4)$$

for all $k \geq M$. Then if $d, d' \in D_m$ for $m > n \geq M$ such that $|d - d'| \leq T2^{-n}$, then

$$\|X(d') - X(d)\| \leq 2 \sum_{j=n+1}^m 2^{-\gamma j}. \quad (30.5)$$

Also, there exists a constant C depending on M such that for all $d, d' \in D$,

$$\|X(d) - X(d')\| \leq C|d - d'|^\gamma.$$

Proof: Suppose $d' < d$. Suppose first $m = n + 1$. Then $d = (k + 1)T2^{-(n+1)}$ and $d' = kT2^{-(n+1)}$. Then from 30.4

$$\|X(d') - X(d)\| \leq 2^{-\gamma(n+1)} \leq 2 \sum_{j=n+1}^{n+1} 2^{-\gamma j}.$$

Suppose the claim 30.5 is true for some $m > n$ and let $d, d' \in D_{m+1}$ with $|d - d'| < T2^{-n}$. If there is no point of D_m between these, then d', d are adjacent points either in D_m or in D_{m+1} . Consequently,

$$\|X(d') - X(d)\| \leq 2^{-\gamma m} < 2 \sum_{j=n+1}^{m+1} 2^{-\gamma j}.$$

Assume therefore, there exist points of D_m between d' and d . Let $d' \leq d'_1 \leq d_1 \leq d$ where d_1, d'_1 are in D_m and d'_1 is the smallest element of D_m which is at least as large as d' and d_1 is the largest element of D_m which is no larger than d . Then $|d' - d'_1| \leq T2^{-(m+1)}$ and $|d_1 - d| \leq T2^{-(m+1)}$ while all of these points are in D_{m+1} which contains D_m . Therefore, from 30.4 and induction,

$$\begin{aligned} & \|X(d') - X(d)\| \\ & \leq \|X(d') - X(d'_1)\| + \|X(d'_1) - X(d_1)\| \\ & \quad + \|X(d_1) - X(d)\| \\ & \leq 2 \times 2^{-\gamma(m+1)} + 2 \sum_{j=n+1}^m 2^{-\gamma j} = 2 \sum_{j=n+1}^{m+1} 2^{-\gamma j} \\ & \leq 2 \left(\frac{2^{-\gamma(n+1)}}{1 - 2^{-\gamma}} \right) = \left(\frac{2T^{-\gamma}}{1 - 2^{-\gamma}} \right) (T2^{-(n+1)})^\gamma \end{aligned} \quad (30.6)$$

It follows the above holds for any $d, d' \in D$ such that $|d - d'| \leq T2^{-n}$ because they are both in some D_m for $m > n$.

Consider the last claim. Let $d, d' \in D, |d - d'| \leq T2^{-M}$. Then d, d' are both in some D_m for $m > M$. The number $|d - d'|$ satisfies

$$T2^{-(n+1)} < |d - d'| \leq T2^{-n}$$

for large enough $n \geq M$. Just pick the first n such that $T2^{-(n+1)} < |d - d'|$. Then from 30.6,

$$\|X(d') - X(d)\| \leq \left(\frac{2T^{-\gamma}}{1 - 2^{-\gamma}} \right) (T2^{-(n+1)})^\gamma \leq \left(\frac{2T^{-\gamma}}{1 - 2^{-\gamma}} \right) (|d - d'|)^\gamma$$

Now $[0, T]$ is covered by 2^M intervals of length $T2^{-M}$ and so for any pair $d, d' \in D$,

$$\|X(d) - X(d')\| \leq C |d - d'|^\gamma$$

where C is a suitable constant depending on 2^M . ■

For $\gamma \leq 1$, you can show, using convexity arguments, that it suffices to have $C = \left(\frac{2T^{-\gamma}}{1 - 2^{-\gamma}} \right)^{1/\gamma} (2^M)^{1-\gamma}$. Of course the case where $\gamma > 1$ is not interesting because it would result in X being a constant.

The following is the amazing Kolmogorov Čentsov continuity theorem [32].

Theorem 30.2.2 Suppose X is a stochastic process on $[0, T]$. Suppose also that there exists a constant, C and positive numbers α, β such that

$$E(\|X(t) - X(s)\|^\alpha) \leq C |t - s|^{1+\beta} \quad (30.7)$$

Then there exists a stochastic process Y such that for a.e. $\omega, t \rightarrow Y(t)(\omega)$ is Hölder continuous with exponent $\gamma < \frac{\beta}{\alpha}$ and for each $t, P(\|X(t) - Y(t)\| > 0) = 0$. (Y is a version of X .)

Proof: Let r_j^m denote $j(\frac{T}{2^m})$ where $j \in \{0, 1, \dots, 2^m\}$. Also let $D_m = \{r_j^m\}_{j=1}^{2^m}$ and $D = \bigcup_{m=1}^{\infty} D_m$. Consider the set,

$$[\|X(t) - X(s)\| > \delta]$$

By 30.7,

$$\begin{aligned} P(\|X(t) - X(s)\| > \delta) \delta^\alpha &\leq \int_{[\|X(t) - X(s)\| > \delta]} \|X(t) - X(s)\|^\alpha dP \\ &\leq C|t - s|^{1+\beta}. \end{aligned} \quad (30.8)$$

Letting $t = r_{j+1}^k, s = r_j^k$, and $\delta = 2^{-\gamma k}$ where $\gamma \in (0, \frac{\beta}{\alpha})$, this yields

$$\begin{aligned} P\left(\|X(r_{j+1}^k) - X(r_j^k)\| > 2^{-\gamma k}\right) &\leq C 2^{\alpha \gamma k} (T 2^{-k})^{1+\beta} \\ &= C T^{1+\beta} 2^{k(\alpha \gamma - (1+\beta))} \end{aligned}$$

There are 2^k of these differences so letting $N_k = \bigcup_{j=1}^{2^k} [\|X(r_{j+1}^k) - X(r_j^k)\| > 2^{-\gamma k}]$ it follows

$$P(N_k) \leq C 2^{\alpha \gamma k} (T 2^{-k})^{1+\beta} 2^k = C 2^{k(\alpha \gamma - \beta)} T^{1+\beta}.$$

Since $\gamma < \beta/\alpha$, $\sum_{k=1}^{\infty} P(N_k) \leq C T^{1+\beta} \sum_{k=1}^{\infty} 2^{k(\alpha \gamma - \beta)} < \infty$ and so by the Borel Cantelli lemma, Lemma 26.1.2, there exists a set of measure zero N , such that if $\omega \notin N$, then ω is in only finitely many N_k . In other words, for $\omega \notin N$, there exists $M(\omega)$ such that if $k \geq M(\omega)$, then for each j ,

$$\|X(r_{j+1}^k)(\omega) - X(r_j^k)(\omega)\| \leq 2^{-\gamma k}. \quad (30.9)$$

It follows from Lemma 30.2.1 that $t \rightarrow X(t)(\omega)$ is Hölder continuous on D with Hölder exponent γ . Note the constant is a measurable function of ω , depending on the number of measurable N_k which contain ω .

By Lemma 30.1.3, one can define $Y(t)(\omega)$ to be the unique function which extends $d \rightarrow X(d)(\omega)$ off D for $\omega \notin N$ and let $Y(t)(\omega) = 0$ if $\omega \in N$. Thus by Lemma 30.1.3 $t \rightarrow Y(t)(\omega)$ is Hölder continuous. Also, $\omega \rightarrow Y(t)(\omega)$ is measurable because it is the pointwise limit of measurable functions

$$Y(t)(\omega) = \lim_{d \rightarrow t} X(d)(\omega) \mathcal{X}_{N^c}(\omega). \quad (30.10)$$

It remains to verify the claim that $Y(t)(\omega) = X(t)(\omega)$ a.e.

$$\mathcal{X}_{[\|Y(t) - X(t)\| > \varepsilon] \cap N^c}(\omega) \leq \liminf_{d \rightarrow t} \mathcal{X}_{[\|X(d) - X(t)\| > \varepsilon] \cap N^c}(\omega)$$

because if $\omega \in N$ both sides are 0 and if $\omega \in N^C$ then the above limit in 30.10 holds and so if $\|Y(t)(\omega) - X(t)(\omega)\| > \varepsilon$, the same is true of $\|X(d)(\omega) - X(t)(\omega)\|$ whenever d is close enough to t and so by Fatou's lemma,

$$\begin{aligned}
 P(\|Y(t) - X(t)\| > \varepsilon) &= \int \mathcal{X}_{\|Y(t) - X(t)\| > \varepsilon \cap N^C}(\omega) dP \\
 &\leq \int \liminf_{d \rightarrow t} \mathcal{X}_{\|X(d) - X(t)\| > \varepsilon}(\omega) dP \\
 &\leq \liminf_{d \rightarrow t} \int \mathcal{X}_{\|X(d) - X(t)\| > \varepsilon}(\omega) dP \\
 &\leq \liminf_{d \rightarrow t} P(\|X(d) - X(t)\|^\alpha > \varepsilon^\alpha) \\
 &\leq \liminf_{d \rightarrow t} \varepsilon^{-\alpha} \int_{\|X(d) - X(t)\|^\alpha > \varepsilon^\alpha} \|X(d) - X(t)\|^\alpha dP \\
 &\leq \liminf_{d \rightarrow t} \frac{C}{\varepsilon^\alpha} |d - t|^{1+\beta} = 0.
 \end{aligned}$$

Therefore,

$$\begin{aligned}
 P(\|Y(t) - X(t)\| > 0) &= P\left(\bigcup_{k=1}^{\infty} \left[\|Y(t) - X(t)\| > \frac{1}{k}\right]\right) \\
 &\leq \sum_{k=1}^{\infty} P\left(\left[\|Y(t) - X(t)\| > \frac{1}{k}\right]\right) = 0. \blacksquare
 \end{aligned}$$

A few observations are interesting. In the proof, the following inequality was obtained.

$$\begin{aligned}
 \|X(d')(\omega) - X(d)(\omega)\| &\leq \frac{2}{T^\gamma(1-2^{-\gamma})} \left(T2^{-(n+1)}\right)^\gamma \\
 &\leq \frac{2}{T^\gamma(1-2^{-\gamma})} (|d - d'|)^\gamma
 \end{aligned}$$

which was so for any $d', d \in D$ with $|d' - d| < T2^{-(M(\omega)+1)}$. Thus the Holder continuous version of X will satisfy

$$\|Y(t)(\omega) - Y(s)(\omega)\| \leq \frac{2}{T^\gamma(1-2^{-\gamma})} (|t - s|)^\gamma$$

provided $|t - s| < T2^{-(M(\omega)+1)}$. Does this translate into an inequality of the form

$$\|Y(t)(\omega) - Y(s)(\omega)\| \leq \frac{2}{T^\gamma(1-2^{-\gamma})} (|t - s|)^\gamma$$

for any pair of points $t, s \in [0, T]$? It seems it does not for any $\gamma < 1$ although it does yield

$$\|Y(t)(\omega) - Y(s)(\omega)\| \leq C(|t - s|)^\gamma$$

where C depends on the number of intervals having length less than $T2^{-(M(\omega)+1)}$ which it takes to cover $[0, T]$. First note that if $\gamma > 1$, then the Holder continuity will imply $t \rightarrow$

$Y(t)(\omega)$ is a constant. Therefore, the only case of interest is $\gamma < 1$. Let s, t be any pair of points and let $s = x_0 < \dots < x_n = t$ where $|x_i - x_{i-1}| < T2^{-(M(\omega)+1)}$. Then

$$\begin{aligned} \|Y(t)(\omega) - Y(s)(\omega)\| &\leq \sum_{i=1}^n \|Y(x_i)(\omega) - Y(x_{i-1})(\omega)\| \\ &\leq \frac{2}{T^\gamma(1-2^{-\gamma})} \sum_{i=1}^n (|x_i - x_{i-1}|)^\gamma \end{aligned} \quad (30.11)$$

How does this compare to $(\sum_{i=1}^n |x_i - x_{i-1}|)^\gamma = |t - s|^\gamma$? This last expression is smaller than the right side of 30.11 for any $\gamma < 1$. Thus for $\gamma < 1$, the constant in the conclusion of the theorem depends on both T and $\omega \notin N$.

In the case where $\alpha \geq 1$, here is another proof of this theorem. It is based on the one in the book by Stroock [56]. This one makes the assumption that $\alpha \geq 1$. It isn't for $\alpha > 0$. This version is sufficient for what is done in this book. The Holder estimate is particularly useful.

Theorem 30.2.3 Suppose X is a stochastic process on $[0, T]$ having values in the Banach space E . Suppose also that there exists a constant C and positive numbers $\alpha, \beta, \alpha \geq 1$, such that

$$E(\|X(t) - X(s)\|^\alpha) \leq C|t - s|^{1+\beta} \quad (30.12)$$

Then there exists a stochastic process Y such that for a.e. $\omega, t \rightarrow Y(t)(\omega)$ is Hölder continuous with exponent $\gamma < \frac{\beta}{\alpha}$ and for each t , $P(\|X(t) - Y(t)\| > 0) = 0$. (Y is a version of X .) Also

$$E\left(\sup_{0 \leq s < t \leq T} \frac{\|Y(t) - Y(s)\|}{(t-s)^\gamma}\right) \leq C$$

where C depends on α, β, T, γ .

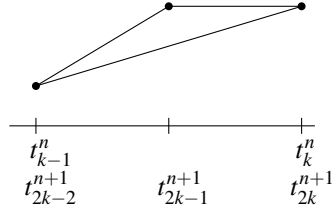
Proof: The proof considers piecewise linear approximations of X which are automatically continuous. These are shown to converge to Y in $L^\alpha(\Omega; C([0, T], E))$ so it will follow that Y must be continuous for a.e. ω . Finally, it is shown that Y is a version of X and is Hölder continuous. In the proof, I will use C to denote a constant which depends on the quantities γ, α, β, T . Let $\{t_k^n\}_{k=0}^{2^n}$ be a uniform partition of the interval $[0, T]$ so that $t_{k+1}^n - t_k^n = T2^{-n}$. Now let

$$M_n \equiv \max_{k \leq 2^n} \|X(t_k^n) - X(t_{k-1}^n)\|$$

Then $M_n^\alpha \leq \max_{k \leq 2^n} \|X(t_k^n) - X(t_{k-1}^n)\|^\alpha \leq \sum_{k=1}^{2^n} \|X(t_k^n) - X(t_{k-1}^n)\|^\alpha$ and so

$$E(M_n^\alpha) \leq \sum_{k=1}^{2^n} C(T2^{-n})^{1+\beta} = C2^n 2^{-n(1+\beta)} = C2^{-n\beta} \quad (30.13)$$

Next denote by X_n the piecewise linear function which results from the values of X at the points t_k^n . Consider the following picture which illustrates a part of the graphs of X_n and X_{n+1} .



Then

$$\begin{aligned} \max_{t \in [0, T]} \|X_{n+1}(t) - X_n(t)\| &\leq \max_{1 \leq k \leq 2^{n+1}} \left\| X(t_{2k-1}^{n+1}) - \frac{X(t_k^n) + X(t_{k-1}^n)}{2} \right\| \\ &\leq \max_{k \leq 2^{n+1}} \left(\frac{1}{2} \|X(t_{2k-1}^{n+1}) - X(t_{2k}^{n+1})\| + \frac{1}{2} \|X(t_{2k-1}^{n+1}) - X(t_{2k-2}^{n+1})\| \right) \leq M_{n+1} \end{aligned}$$

Denote by $\|\cdot\|_\infty$ the usual norm in $C([0, T], E)$, $\max_{t \in [0, T]} \|Z(t)\| \equiv \|Z\|_\infty$. Then from what was just established,

$$E(\|X_{n+1} - X_n\|_\infty^\alpha) = \int_\Omega \|X_{n+1} - X_n\|_\infty^\alpha dP \leq E(M_{n+1}^\alpha) = C 2^{-n\beta}$$

which shows that

$$\|X_{n+1} - X_n\|_{L^\alpha(\Omega; C([0, T], E))} = \left(\int_\Omega \|X_{n+1} - X_n\|_\infty^\alpha dP \right)^{1/\alpha} \leq C \left(2^{(\beta/\alpha)} \right)^{-n}$$

Since $\alpha \geq 1$, we can use the triangle inequality and conclude

$$\begin{aligned} \|X_m - X_n\|_{L^\alpha(\Omega; C([0, T], E))} &\leq \\ \sum_{k=n}^{\infty} C \left(2^{(\beta/\alpha)} \right)^{-k} &\leq C \frac{\left(2^{(\beta/\alpha)} \right)^{-n}}{1 - 2^{-(\beta/\alpha)}} = C \left(2^{(\beta/\alpha)} \right)^{-n} \end{aligned} \quad (30.14)$$

Thus $\{X_n\}$ is a Cauchy sequence in $L^\alpha(\Omega; C([0, T], E))$ and so it converges to some Y in this space, a subsequence converging pointwise. Then from Fatou's lemma,

$$\|Y - X_n\|_{L^\alpha(\Omega; C([0, T], E))} \leq C \left(2^{(\beta/\alpha)} \right)^{-n}. \quad (30.15)$$

Also, for a.e. ω , $t \rightarrow Y(t)$ is in $C([0, T], E)$. It remains to verify that $Y(t) = X(t)$ a.e.

From the construction, it follows that for any n and $m \geq n$, $Y(t_k^n) = X_m(t_k^n) = X(t_k^n)$. Thus

$$\begin{aligned} \|Y(t) - X(t)\| &\leq \|Y(t) - Y(t_k^n)\| + \|Y(t_k^n) - X(t)\| \\ &= \|Y(t) - Y(t_k^n)\| + \|X(t_k^n) - X(t)\| \end{aligned}$$

Now from the hypotheses of the theorem,

$$P(\|X(t_k^n) - X(t)\|^\alpha > \varepsilon) \leq \frac{1}{\varepsilon} E(\|X(t_k^n) - X(t)\|^\alpha) \leq \frac{C}{\varepsilon} |t_k^n - t|^{1+\beta}$$

Thus, there exists a sequence of mesh points $\{s_n\}$ converging to t such that

$$P(\|X(s_n) - X(t)\|^\alpha > 2^{-n}) \leq 2^{-n}$$

Then by the Borel Cantelli lemma, there is a set of measure zero N such that for $\omega \notin N$, $\|X(s_n) - X(t)\|^\alpha \leq 2^{-n}$ for all n large enough. Then

$$\|Y(t) - X(t)\| \leq \|Y(t) - Y(s_n)\| + \|X(s_n) - X(t)\|$$

which shows that, by continuity of Y , for ω not in an exceptional set of measure zero, $\|Y(t) - X(t)\| = 0$.

It remains to verify the assertion about Hölder continuity of Y . Let $0 \leq s < t \leq T$. Then for some n ,

$$2^{-(n+1)}T \leq t - s \leq 2^{-n}T \quad (30.16)$$

Thus

$$\begin{aligned} \|Y(t) - Y(s)\| &\leq \|Y(t) - X_n(t)\| + \|X_n(t) - X_n(s)\| + \|X_n(s) - Y(s)\| \\ &\leq 2 \sup_{\tau \in [0, T]} \|Y(\tau) - X_n(\tau)\| + \|X_n(t) - X_n(s)\| \end{aligned} \quad (30.17)$$

Now

$$\frac{\|X_n(t) - X_n(s)\|}{t - s} \leq \frac{\|X_n(t) - X_n(s)\|}{2^{-(n+1)}T}$$

From 30.16 a picture like the following must hold in which $t_{k-1}^{n+1} \leq s < t_k^{n+1} < t \leq t_{k+1}^{n+1}$.

$$\begin{array}{ccccccc} | & & | & & | & & | \\ t_{k-1}^{n+1} & s & t_k^{n+1} & t & t_{k+1}^{n+1} \end{array}$$

Therefore, from the above, 30.16,

$$\begin{aligned} \frac{\|X_n(t) - X_n(s)\|}{t - s} &\leq \frac{\|X(t_{k-1}^{n+1}) - X(t_k^{n+1})\| + \|X(t_k^{n+1}) - X(t_{k+1}^{n+1})\|}{2^{-(n+1)}T} \\ &\leq C2^n M_{n+1} \end{aligned}$$

It follows from 30.17,

$$\|Y(t) - Y(s)\| \leq 2\|Y - X_n\|_\infty + C2^n M_{n+1}(t - s)$$

Next, letting $\gamma < \beta/\alpha$, and using 30.16,

$$\begin{aligned} \frac{\|Y(t) - Y(s)\|}{(t - s)^\gamma} &\leq 2(T^{-1}2^{n+1})^\gamma \|Y - X_n\|_\infty + C2^n (2^{-n})^{1-\gamma} M_{n+1} \\ &= C2^{n\gamma} (\|Y - X_n\|_\infty + M_{n+1}) \end{aligned}$$

The above holds for any s, t satisfying 30.16. Then

$$\sup \left\{ \frac{\|Y(t) - Y(s)\|}{(t - s)^\gamma}, 0 \leq s < t \leq T, |t - s| \in [2^{-(n+1)}T, 2^{-n}T] \right\}$$

$$\leq C2^{n\gamma}(\|Y - X_n\|_\infty + M_{n+1})$$

Denote by P_n the ordered pairs (s, t) satisfying the above condition that

$$0 \leq s < t \leq T, |t - s| \in \left[2^{-(n+1)}T, 2^{-n}T\right]$$

and also

$$\sup_{(s,t) \in P_n} \frac{\|Y(t) - Y(s)\|}{(t-s)^\gamma} \leq C2^{n\gamma}(\|Y - X_n\|_\infty + M_{n+1})$$

Note that the union of the P_n pertains to all (s, t) with $|t - s| \leq T/2$. If $|t - s| > T/2$, then $E\left(\frac{\|Y(t) - Y(s)\|}{|t-s|^\gamma}\right) \leq \left(\frac{2}{T}\right)^\gamma 2\|Y\|_{L^1(\Omega; C([0, T]; E))}$ so the desired condition holds and we can ignore this case.

Thus for *a.e.* ω , and for all n ,

$$\left(\sup_{(s,t) \in P_n} \frac{\|Y(t) - Y(s)\|}{(t-s)^\gamma}\right)^\alpha \leq C \sum_{k=0}^{\infty} 2^{k\alpha\gamma} (\|Y - X_k\|_\infty^\alpha + M_{k+1}^\alpha)$$

Note that n is arbitrary. Hence

$$\begin{aligned} \sup_{0 \leq s < t \leq T} \left(\frac{\|Y(t) - Y(s)\|}{(t-s)^\gamma}\right)^\alpha &\leq \\ \sup_n \sup_{(s,t) \in P_n} \left(\frac{\|Y(t) - Y(s)\|}{(t-s)^\gamma}\right)^\alpha &\leq \sup_n \left(\sup_{(s,t) \in P_n} \frac{\|Y(t) - Y(s)\|}{(t-s)^\gamma}\right)^\alpha \\ &\leq \sum_{k=0}^{\infty} C2^{k\alpha\gamma} (\|Y - X_k\|_\infty^\alpha + M_{k+1}^\alpha) \end{aligned}$$

By continuity of Y , the result on the left is unchanged if the ordered pairs are restricted to lie in $\mathbb{Q} \cap [0, T] \times \mathbb{Q} \cap [0, T]$, a countable set. Thus the left side is measurable. It follows from 30.13 and 30.15 which say

$$\|Y - X_k\|_{L^\alpha(\Omega; C([0, T], E))} \leq C \left(2^{(\beta/\alpha)}\right)^{-k}, E(M_k^\alpha) \leq C2^{-k\beta}$$

that

$$E\left(\sup_{0 \leq s < t \leq T} \left(\frac{\|Y(t) - Y(s)\|}{(t-s)^\gamma}\right)^\alpha\right) \leq \sum_{k=0}^{\infty} C2^{k\alpha\gamma} 2^{-\beta k} \equiv C < \infty$$

because $\alpha\gamma - \beta < 0$. By continuity of Y , there are no measurability concerns in taking the above expectation. Note that this implies, since $\alpha \geq 1$,

$$\begin{aligned} E\left(\sup_{0 \leq s < t \leq T} \frac{\|Y(t) - Y(s)\|}{(t-s)^\gamma}\right) &\leq \left(E\left(\sup_{0 \leq s < t \leq T} \left(\frac{\|Y(t) - Y(s)\|}{(t-s)^\gamma}\right)^\alpha\right)\right)^{1/\alpha} \\ &\leq C^{1/\alpha} \equiv C \end{aligned}$$

Now

$$P\left(\sup_{0 \leq s < t \leq T} \left(\frac{\|Y(t) - Y(s)\|}{(t-s)^\gamma}\right)^\alpha > 2^{\alpha k}\right) \leq \frac{1}{2^{\alpha k}} C$$

and so there exists a set of measure zero N such that for $\omega \notin N$,

$$\sup_{0 \leq s < t \leq T} \left(\frac{\|Y(t) - Y(s)\|}{(t-s)^\gamma} \right)^\alpha \leq 2^{\alpha k}$$

for all k large enough. Pick such a k , depending on $\omega \notin N$. Then for any s, t ,

$$\frac{\|Y(t) - Y(s)\|}{(t-s)^\gamma} \leq 2^k$$

and so, this has shown that for $\omega \notin N$, $\|Y(t) - Y(s)\| \leq C(\omega)(t-s)^\gamma$ ■

Note that if $X(t)$ is known to be continuous off a set of measure zero, then the piecewise linear approximations converge to $X(t)$ in $C([0, T], E)$ off this set of measure zero. Therefore, it must be that off a set of measure zero, $Y(t) = X(t)$ and so in fact $X(t)$ is Holder continuous off a set of measure zero and the condition on expectation also must hold, that is

$$E \left(\sup_{0 \leq s < t \leq T} \frac{\|X(t) - X(s)\|}{(t-s)^\gamma} \right) \leq C.$$

30.3 Filtrations

Instead of having a sequence of σ algebras, one can consider an increasing collection of σ algebras indexed by $t \in \mathbb{R}$. This is called a filtration.

Definition 30.3.1 Let X be a stochastic process defined on an interval, $I = [0, T]$ or $[0, \infty)$. Suppose the probability space, (Ω, \mathcal{F}, P) has an increasing family of σ algebras, $\{\mathcal{F}_t\}$. This is called a filtration. If for arbitrary $t \in I$ the random variable $X(t)$ is \mathcal{F}_t measurable, then X is said to be adapted to the filtration $\{\mathcal{F}_t\}$. Denote by \mathcal{F}_{t+} the intersection of all \mathcal{F}_s for $s > t$. The filtration is normal if

1. \mathcal{F}_0 contains all $A \in \mathcal{F}$ such that $P(A) = 0$
2. $\mathcal{F}_t = \mathcal{F}_{t+}$ for all $t \in I$.

X is called progressively measurable if for every $t \in I$, the mapping

$$(s, \omega) \in [0, t] \times \Omega, (s, \omega) \rightarrow X(s, \omega)$$

is $B([0, t]) \times \mathcal{F}_t$ measurable.

Thus X is progressively measurable means

$$(s, \omega) \rightarrow \mathcal{X}_{[0, t]}(s) X(s, \omega)$$

is $B([0, t]) \times \mathcal{F}_t$ measurable. As an example of a normal filtration, here is an example.

Example 30.3.2 For example, you could have a stochastic process, $X(t)$ and you could define

$$\mathcal{G}_t \equiv \overline{\sigma(X(s) : s \leq t)},$$

the completion of the smallest σ algebra such that each $X(s)$ is measurable for all $s \leq t$. This gives an example of a filtration to which $X(t)$ is adapted which satisfies 1. More generally, suppose $X(t)$ is adapted to a filtration, \mathcal{G}_t . Define

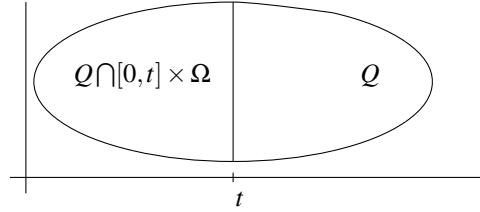
$$\mathcal{F}_t \equiv \cap_{s>t} \mathcal{G}_s$$

Then

$$\mathcal{F}_{t+} \equiv \cap_{s>t} \mathcal{F}_s = \cap_{s>t} \cap_{r>s} \mathcal{G}_r = \cap_{s>t} \mathcal{F}_s \equiv \mathcal{F}_t.$$

and each $X(t)$ is measurable with respect to \mathcal{F}_t . Thus there is no harm in assuming a stochastic process adapted to a filtration can be modified so the filtration is normal. Also note that \mathcal{F}_t defined this way will be complete so if $A \in \mathcal{F}_t$ has $P(A) = 0$ and if $B \subseteq A$, then $B \in \mathcal{F}_t$ also. This is because this relation between the sets and the probability of A being zero, holds for this pair of sets when considered as elements of each \mathcal{G}_s for $s > t$. Hence $B \in \mathcal{G}_s$ for each $s > t$ and is therefore one of the sets in \mathcal{F}_t .

What is the description of a progressively measurable set?



It means that for Q progressively measurable, $Q \cap [0, t] \times \Omega$ as shown in the above picture is $\mathcal{B}([0, t]) \times \mathcal{F}_t$ measurable. It is like saying a little more descriptively that the function is progressively product measurable.

I shall generally assume the filtration is normal.

Observation 30.3.3 *If X is progressively measurable, then it is adapted. Furthermore the progressively measurable sets, those $E \cap [0, T] \times \Omega$ for which \mathcal{X}_E is progressively measurable form a σ algebra.*

To see why this is, consider X progressively measurable and fix t . Then $(s, \omega) \rightarrow X(s, \omega)$ for $(s, \omega) \in [0, t] \times \Omega$ is given to be $\mathcal{B}([0, t]) \times \mathcal{F}_t$ measurable, the ordinary product measure and so fixing any $s \in [0, t]$, it follows the resulting function of ω is \mathcal{F}_t measurable. In particular, this is true upon fixing $s = t$. Thus $\omega \rightarrow X(t, \omega)$ is \mathcal{F}_t measurable and so $X(t)$ is adapted.

A set $E \subseteq [0, T] \times \Omega$ is progressively measurable means that \mathcal{X}_E is progressively measurable. That is \mathcal{X}_E restricted to $[0, t] \times \Omega$ is $\mathcal{B}([0, t]) \times \mathcal{F}_t$ measurable. In other words, E is progressively measurable if

$$E \cap ([0, t] \times \Omega) \in \mathcal{B}([0, t]) \times \mathcal{F}_t.$$

If E_i is progressively measurable, does it follow that $E \equiv \cup_{i=1}^{\infty} E_i$ is also progressively measurable? Yes.

$$E \cap ([0, t] \times \Omega) = \cup_{i=1}^{\infty} E_i \cap ([0, t] \times \Omega) \in \mathcal{B}([0, t]) \times \mathcal{F}_t$$

because each set in the union is in $\mathcal{B}([0, t]) \times \mathcal{F}_t$. If E is progressively measurable, is E^C ?

$$E^C \cap ([0, t] \times \Omega) \cup \overbrace{(E \cap ([0, t] \times \Omega))}^{\in \mathcal{B}([0, t]) \times \mathcal{F}_t} = \overbrace{[0, t] \times \Omega}^{\in \mathcal{B}([0, t]) \times \mathcal{F}_t}$$

and so $E^C \cap ([0, t] \times \Omega) \in \mathcal{B}([0, t]) \times \mathcal{F}_t$. Thus the progressively measurable sets are a σ algebra.

Another observation of interest is in the following lemma.

Lemma 30.3.4 *Suppose Q is in $\mathcal{B}([0, a]) \times \mathcal{F}_r$. Then if $b \geq a$ and $t \geq r$, then Q is also in $\mathcal{B}([0, b]) \times \mathcal{F}_t$.*

Proof: Consider a measurable rectangle $A \times B$ where $A \in \mathcal{B}([0, a])$ and $B \in \mathcal{F}_r$. Is it true that $A \times B \in \mathcal{B}([0, b]) \times \mathcal{F}_t$? This reduces to the question whether $A \in \mathcal{B}([0, b])$. If A is an interval, it is clear that $A \in \mathcal{B}([0, b])$. Consider the π system of intervals and let \mathcal{G} denote those Borel sets $A \in \mathcal{B}([0, a])$ such that $A \in \mathcal{B}([0, b])$. If $A \in \mathcal{G}$, then $[0, b] \setminus A \in \mathcal{B}([0, b])$ by assumption (the difference of Borel sets is surely Borel). However, this set equals

$$([0, a] \setminus A) \cup (a, b]$$

and so

$$[0, b] = ([0, a] \setminus A) \cup (a, b] \cup A$$

The set on the left is in $\mathcal{B}([0, b])$ and the sets on the right are disjoint and two of them are also in $\mathcal{B}([0, b])$. Therefore, the third, $([0, a] \setminus A)$ is in $\mathcal{B}([0, b])$. It is obvious that \mathcal{G} is closed with respect to countable disjoint unions. Therefore, by Lemma 9.3.2, Dynkin's lemma, $\mathcal{G} \supseteq \sigma(\text{Intervals}) = \mathcal{B}([0, a])$.

Therefore, such a measurable rectangle $A \times B$ where $A \in \mathcal{B}([0, a])$ and $B \in \mathcal{F}_r$ is in $\mathcal{B}([0, b]) \times \mathcal{F}_t$ and in fact it is a measurable rectangle in $\mathcal{B}([0, b]) \times \mathcal{F}_t$. Now let \mathcal{H} denote all these measurable rectangles $A \times B$ where $A \in \mathcal{B}([0, a])$ and $B \in \mathcal{F}_r$. Let \mathcal{G} (new \mathcal{G}) denote those sets Q of $\mathcal{B}([0, a]) \times \mathcal{F}_r$ which are in $\mathcal{B}([0, b]) \times \mathcal{F}_t$. Then if $Q \in \mathcal{G}$,

$$Q \cup ([0, a] \times \Omega \setminus Q) \cup (a, b] \times \Omega = [a, b] \times \Omega$$

Then the sets are disjoint and all but $[0, a] \times \Omega \setminus Q$ are in $\mathcal{B}([0, b]) \times \mathcal{F}_t$. Therefore, this one is also in $\mathcal{B}([0, b]) \times \mathcal{F}_t$. If $Q_i \in \mathcal{G}$ and the Q_i are disjoint, then $\cup_i Q_i$ is also in $\mathcal{B}([0, b]) \times \mathcal{F}_t$ and so \mathcal{G} is closed with respect to countable disjoint unions and complements. Hence $\mathcal{G} \supseteq \sigma(\mathcal{H}) = \mathcal{B}([0, a]) \times \mathcal{F}_r$ which shows

$$\mathcal{B}([0, a]) \times \mathcal{F}_r \subseteq \mathcal{B}([0, b]) \times \mathcal{F}_t \quad \blacksquare$$

A significant observation is the following which states that the integral of a progressively measurable function is progressively measurable.

Proposition 30.3.5 *Suppose $X : [0, T] \times \Omega \rightarrow E$ where E is a separable Banach space. Also suppose that $X(\cdot, \omega) \in L^1([0, T], E)$ for each ω . Here \mathcal{F}_t is a filtration and with respect to this filtration, X is progressively measurable. Then*

$$(t, \omega) \rightarrow \int_0^t X(s, \omega) ds$$

is also progressively measurable.

Proof: Suppose $Q \in [0, T] \times \Omega$ is progressively measurable. This means for each t ,

$$Q \cap [0, t] \times \Omega \in \mathcal{B}([0, t]) \times \mathcal{F}_t$$

What about $(s, \omega) \in [0, t] \times \Omega$, $(s, \omega) \rightarrow \int_0^s \mathcal{X}_Q dr$? Is that function on the right $\mathcal{B}([0, t]) \times \mathcal{F}_t$ measurable? We know that $Q \cap [0, s] \times \Omega$ is $\mathcal{B}([0, s]) \times \mathcal{F}_s$ measurable and hence $\mathcal{B}([0, t]) \times \mathcal{F}_t$ measurable. When you integrate a product measurable function, you do get one which is product measurable. Therefore, this function must be $\mathcal{B}([0, t]) \times \mathcal{F}_t$ measurable. This shows that $(t, \omega) \rightarrow \int_0^t \mathcal{X}_Q(s, \omega) ds$ is progressively measurable. Here is a claim which was just used.

Claim: If Q is $\mathcal{B}([0, t]) \times \mathcal{F}_t$ measurable, then $(s, \omega) \rightarrow \int_0^s \mathcal{X}_Q dr$ is also $\mathcal{B}([0, t]) \times \mathcal{F}_t$ measurable.

Proof of claim: First consider $A \times B$ where $A \in \mathcal{B}([0, t])$ and $B \in \mathcal{F}_t$. Then

$$\int_0^s \mathcal{X}_{A \times B} dr = \int_0^s \mathcal{X}_A \mathcal{X}_B dr = \mathcal{X}_B(\omega) \int_0^s \mathcal{X}_A(s) dr$$

This is clearly $\mathcal{B}([0, t]) \times \mathcal{F}_t$ measurable. It is the product of a continuous function of s with the indicator function of a set in \mathcal{F}_t . Now let

$$\mathcal{G} \equiv \left\{ Q \in \mathcal{B}([0, t]) \times \mathcal{F}_t : (s, \omega) \rightarrow \int_0^s \mathcal{X}_Q(r, \omega) dr \text{ is } \mathcal{B}([0, t]) \times \mathcal{F}_t \text{ measurable} \right\}$$

Then it was just shown that \mathcal{G} contains the measurable rectangles. It is also clear that \mathcal{G} is closed with respect to countable disjoint unions and complements. Therefore, $\mathcal{G} \supseteq \sigma(\mathcal{K}_t) = \mathcal{B}([0, t]) \times \mathcal{F}_t$ where \mathcal{K}_t denotes the measurable rectangles $A \times B$ where $B \in \mathcal{F}_t$ and $A \in \mathcal{B}([0, t]) = \mathcal{B}([0, T]) \cap [0, t]$. This proves the claim.

Thus if Q is progressively measurable, $(s, \omega) \rightarrow \int_0^s \mathcal{X}_Q(r, \omega) dr \equiv f(s, \omega)$ is progressively measurable because for $(s, \omega) \in [0, t] \times \Omega$, $(s, \omega) \rightarrow f(s, \omega)$ is $\mathcal{B}([0, t]) \times \mathcal{F}_t$ measurable. This is what was to be proved in this special case.

Now consider the conclusion of the proposition. By considering the positive and negative parts of $\phi(X)$ for $\phi \in E'$, and using Pettis theorem, it suffices to consider the case where $X \geq 0$. Then there exists an increasing sequence of progressively measurable simple functions $\{X_n\}$ converging pointwise to X . From what was just shown,

$$(t, \omega) \rightarrow \int_0^t X_n ds$$

is progressively measurable. Hence, by the monotone convergence theorem, $(t, \omega) \rightarrow \int_0^t X ds$ is also progressively measurable. ■

What else can you do to something which is progressively measurable and obtain something which is progressively measurable? It turns out that shifts in time can preserve progressive measurability. Let \mathcal{F}_t be a filtration on $[0, T]$ and extend the filtration to be equal to \mathcal{F}_0 and \mathcal{F}_T for $t < 0$ and $t > T$, respectively. Recall the following definition of progressively measurable sets.

Definition 30.3.6 Denote by \mathcal{P} those sets Q in $\mathcal{F}_T \times \mathcal{B}([0, T])$ such that for $t \in [-\infty, T]$

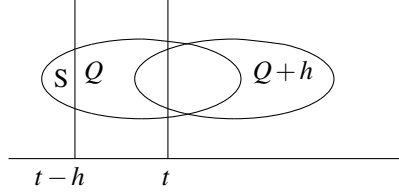
$$\Omega \times (-\infty, t] \cap Q \in \mathcal{F}_t \times \mathcal{B}((-\infty, t]).$$

Lemma 30.3.7 Define $Q+h$ as

$$Q+h \equiv \{(t+h, \omega) : (t, \omega) \in Q\}.$$

Then if $Q \in \mathcal{P}$, it follows that $Q+h \in \mathcal{P}$.

Proof: This is most easily seen through the use of the following diagram. In this diagram, Q is in \mathcal{P} so it is progressively measurable.



By definition, S in the picture is $\mathcal{B}((-\infty, t-h]) \times \mathcal{F}_{t-h}$ measurable. Hence $S+h \equiv Q+h \cap \Omega \times (-\infty, t]$ is $\mathcal{B}((-\infty, t]) \times \mathcal{F}_{t-h}$ measurable. To see this, note that if $B \times A \in \mathcal{B}((-\infty, t-h]) \times \mathcal{F}_{t-h}$, then translating it by h gives a set in $\mathcal{B}((-\infty, t]) \times \mathcal{F}_{t-h}$. Then if \mathcal{G} consists of sets S in $\mathcal{B}((-\infty, t-h]) \times \mathcal{F}_{t-h}$ for which $S+h$ is in $\mathcal{B}((-\infty, t]) \times \mathcal{F}_{t-h}$, \mathcal{G} is closed with respect to countable disjoint unions and complements. Thus, \mathcal{G} equals $\mathcal{B}((-\infty, t-h]) \times \mathcal{F}_{t-h}$. In particular, it contains the set S just described. ■

Now for $h > 0$,

$$\tau_h f(t) \equiv \begin{cases} f(t-h) & \text{if } t \geq h, \\ 0 & \text{if } t < h. \end{cases}$$

Lemma 30.3.8 Let $Q \in \mathcal{P}$. Then $\tau_h \mathcal{X}_Q$ is \mathcal{P} measurable.

Proof: If $\tau_h \mathcal{X}_Q(t, \omega) = 1$, then you need to have $(t-h, \omega) \in Q$ and so $(t, \omega) \in Q+h$. Thus

$$\tau_h \mathcal{X}_Q = \mathcal{X}_{Q+h},$$

which is \mathcal{P} measurable since $Q \in \mathcal{P}$. In general,

$$\tau_h \mathcal{X}_Q = \mathcal{X}_{[h, T] \times \Omega} \mathcal{X}_{Q+h},$$

which is \mathcal{P} measurable. ■

This lemma implies the following.

Lemma 30.3.9 Let $f(t, \omega)$ have values in a separable Banach space and suppose f is \mathcal{P} measurable. Then $\tau_h f$ is \mathcal{P} measurable.

Proof: Taking values in a separable Banach space and being \mathcal{P} measurable, f is the pointwise limit of \mathcal{P} measurable simple functions. If s_n is one of these, then from the above lemmas, $\tau_h s_n$ is \mathcal{P} measurable. Then, letting $n \rightarrow \infty$, it follows that $\tau_h f$ is \mathcal{P} measurable. ■

The following is similar to Proposition 30.1.2. It shows that under pretty weak conditions, an adapted process has a progressively measurable adapted version.

Proposition 30.3.10 Let X be a stochastically continuous adapted process for a normal filtration defined on a closed interval, $I \equiv [0, T]$. Then X has a progressively measurable adapted version.

Proof: By Lemma 30.1.1 X is uniformly stochastically continuous and so there exists a sequence of positive numbers, $\{\rho_n\}$ such that if $|s - t| < \rho_n$, then

$$P\left(\left[\|X(t) - X(s)\| \geq \frac{1}{2^n}\right]\right) \leq \frac{1}{2^n}. \quad (30.18)$$

Then let $\{t_0^n, t_1^n, \dots, t_{m_n}^n\}$ be a partition of $[0, T]$ in which $|t_i^n - t_{i-1}^n| < \rho_n$. Now define X_n as follows:

$$\begin{aligned} X_n(t)(\omega) &\equiv \sum_{i=1}^{m_n} X(t_{i-1}^n)(\omega) \mathcal{X}_{[t_{i-1}^n, t_i^n)}(t) \\ X_n(T) &\equiv X(T). \end{aligned}$$

Then $(s, \omega) \rightarrow X_n(s, \omega)$ for $(s, \omega) \in [0, t] \times \Omega$ is obviously $B([0, t]) \times \mathcal{F}_t$ measurable. Consider the set, A on which $\{X_n(t, \omega)\}$ is a Cauchy sequence. This set is of the form

$$A = \cap_{n=1}^{\infty} \cup_{m=1}^{\infty} \cap_{p, q \geq m} \left[\|X_p - X_q\| < \frac{1}{n} \right]$$

and so it is a $B(I) \times \mathcal{F}$ measurable set and $A \cap [0, t] \times \Omega$ is $B([0, t]) \times \mathcal{F}_t$ measurable for each $t \leq T$ because each X_q in the above has the property that its restriction to $[0, t] \times \Omega$ is $B([0, t]) \times \mathcal{F}_t$ measurable. Now define

$$Y(t, \omega) \equiv \begin{cases} \lim_{n \rightarrow \infty} X_n(t, \omega) & \text{if } (t, \omega) \in A \\ 0 & \text{if } (t, \omega) \notin A \end{cases}$$

I claim that for each t , $Y(t, \omega) = X(t, \omega)$ for a.e. ω . To see this, consider 30.18. From the construction of X_n , it follows that for each t ,

$$P\left(\left[\|X_n(t) - X(t)\| \geq \frac{1}{2^n}\right]\right) \leq \frac{1}{2^n}$$

Also, for a fixed t , if $X_n(t, \omega)$ fails to converge to $X(t, \omega)$, then ω must be in infinitely many of the sets,

$$B_n \equiv \left[\|X_n(t) - X(t)\| \geq \frac{1}{2^n} \right]$$

which is a set of measure zero by the Borel Cantelli lemma. Recall why this is so.

$$P(\cap_{k=1}^{\infty} \cup_{n=k}^{\infty} B_n) \leq \sum_{n=k}^{\infty} P(B_n) < \frac{1}{2^{k-1}}$$

Therefore, for each t , $(t, \omega) \in A$ for a.e. ω . Hence $X(t) = Y(t)$ a.e. and so Y is a measurable version of X . Y is adapted because the filtration is normal and hence \mathcal{F}_t contains all sets of measure zero. Therefore, $Y(t)$ differs from $X(t)$ on a set which is \mathcal{F}_t measurable. ■

There is a more specialized situation in which the measurability of a stochastic process automatically implies it is adapted. Furthermore, this can be defined easily in terms of a π system of sets.

Definition 30.3.11 Let \mathcal{F}_t be a filtration on (Ω, \mathcal{F}, P) and denote by \mathcal{P}_{∞} the smallest σ algebra of sets of $[0, \infty) \times \Omega$ containing the sets

$$(s, t] \times F, F \in \mathcal{F}_s \quad \{0\} \times F, F \in \mathcal{F}_0.$$

This is a lot like product measure except one of the σ algebras is changing.

Proof: Let $s_0 > 0$ and define

where

The diagram shows a coordinate system with a horizontal axis labeled Ω and a vertical axis. A horizontal line is drawn at a height labeled s_0 on the vertical axis. This line intersects an ellipse. The portion of this horizontal line that lies within the ellipse is highlighted with a red segment. Dotted vertical lines extend from the endpoints of this red segment down to the horizontal axis, where the interval between them is labeled S_{s_0} .

It is clear \mathcal{G}_{s_0} is a σ algebra. The next step is to show \mathcal{G}_{s_0} contains the sets

and

It is clear $\{0\} \times F$ is contained in \mathcal{G}_{s_0} because $(\{0\} \times F)_{s_0} = \emptyset \in \mathcal{F}_{s_0}$. Similarly, if $s \geq s_0$ or if $s, t < s_0$ then $((s, t] \times F)_{s_0} = \emptyset \in \mathcal{F}_{s_0}$. The only case left is for $s < s_0$ and $t \geq s_0$. In this case, letting $A_s \in \mathcal{F}_s$, $((s, t] \times A_s)_{s_0} = A_s \in \mathcal{F}_s \subseteq \mathcal{F}_{s_0}$. Therefore, \mathcal{G}_{s_0} contains all the sets of the form given in 30.19 and 30.20 and so since \mathcal{P}_∞ is the smallest σ algebra containing these sets, it follows $\mathcal{P}_\infty = \mathcal{G}_{s_0}$. The case where $s_0 = 0$ is entirely similar but shorter.

Therefore, if X is predictable, letting $A \in \mathcal{B}(E)$, $X^{-1}(A) \in \mathcal{P}_\infty$ or \mathcal{P}_T and so

showing $X(t)$ is \mathcal{F}_t adapted. ■

Proposition 30.3.13 *Let \mathcal{P} denote the predictable σ algebra and let \mathcal{R} denote the progressively measurable σ algebra. Then $\mathcal{P} \subset \mathcal{R}$.*

Proof: Let \mathcal{G} denote those sets of \mathcal{P} such that they are also in \mathcal{R} . Then \mathcal{G} clearly contains the π system of sets $\{0\} \times A, A \in \mathcal{F}_0$, and $(s, t] \times A, A \in \mathcal{F}_s$. Furthermore, \mathcal{G} is closed with respect to countable disjoint unions and complements. It follows \mathcal{G} contains the σ algebra generated by this π systems which is \mathcal{P} . ■

Proposition 30.3.14 *Let $X(t)$ be a stochastic process having values in E a complete metric space and let it be \mathcal{F}_t adapted and left continuous. Then it is predictable. Also, if $X(t)$ is stochastically continuous and adapted on $[0, T]$, then it has a predictable version.*

Proof: Define $I_{m,k} \equiv ((k-1)2^{-m}T, k2^{-m}T]$ if $k \geq 1$ and $I_{m,0} = \{0\}$ if $k = 1$. Then define

$$\begin{aligned} X_m(t) &\equiv \sum_{k=1}^{2^m} X((k-1)2^{-m}T) \mathcal{I}_{((k-1)2^{-m}T, k2^{-m}T]}(t) \\ &\quad + X(0) \mathcal{I}_{[0,0]}(t) \end{aligned}$$

Here the sum means that $X_m(t)$ has value $X((k-1)2^{-m}T)$ on the interval

$$((k-1)2^{-m}T, k2^{-m}T].$$

Thus X_m is predictable because each term in the sum is. Thus

$$\begin{aligned} X_m^{-1}(U) &= \cup_{k=1}^{2^m} (X((k-1)2^{-m}T) \mathcal{I}_{((k-1)2^{-m}T, k2^{-m}T]}^{-1}(U)) \\ &= \cup_{k=1}^{2^m} ((k-1)2^{-m}T, k2^{-m}T] \times (X((k-1)2^{-m}T))^{-1}(U), \end{aligned}$$

a finite union of predictable sets. Since X is left continuous,

$$X(t, \omega) = \lim_{m \rightarrow \infty} X_m(t, \omega)$$

and so X is predictable.

Next consider the other claim. Since X is stochastically continuous on $[0, T]$, it is uniformly stochastically continuous on this interval by Lemma 30.1.1. Therefore, there exists a sequence of partitions of $[0, T]$, the m^{th} being

$$0 = t_{m,0} < t_{m,1} < \cdots < t_{m,n(m)} = T$$

such that for X_m defined as above, then for each t

$$P([d(X_m(t), X(t)) \geq 2^{-m}]) \leq 2^{-m} \quad (30.21)$$

Then as above, X_m is predictable. Let A denote those points of \mathcal{P}_T at which $X_m(t, \omega)$ converges. Thus A is a predictable set because it is just the set where $X_m(t, \omega)$ is a Cauchy sequence. Now define the predictable function Y

$$Y(t, \omega) \equiv \begin{cases} \lim_{m \rightarrow \infty} X_m(t, \omega) & \text{if } (t, \omega) \in A \\ 0 & \text{if } (t, \omega) \notin A \end{cases}$$

From 30.21 it follows from the Borel Cantelli lemma that for fixed t , the set of ω which are in infinitely many of the sets,

$$[d(X_m(t), X(t)) \geq 2^{-m}]$$

has measure zero. Therefore, for each t , there exists a set of measure zero, $N(t)$ such that for $\omega \notin N(t)$ and all m large enough

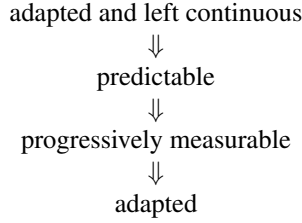
$$d(X_m(t, \omega), X(t, \omega)) < 2^{-m}$$

Hence for $\omega \notin N(t)$, $(t, \omega) \in A$ and so $X_m(t, \omega) \rightarrow Y(t, \omega)$ which shows

$$d(Y(t, \omega), X(t, \omega)) = 0 \text{ if } \omega \notin N(t).$$

The predictable version of $X(t)$ is $Y(t)$. ■

Here is a summary of what has been shown above.



Also

stochastically continuous and adapted \implies progressively measurable version

30.4 Martingales and Sub-Martingales

This was done earlier for discreet martingales. The idea here is to consider indiscreet (What a word to use for a martingale!) ones.

Definition 30.4.1 Let X be a stochastic process defined on an interval I with values in a separable Banach space, E . It is called integrable if $E(\|X(t)\|) < \infty$ for each $t \in I$. Also let \mathcal{F}_t be a filtration. An integrable and adapted stochastic process X is called a martingale if for $s \leq t$

$$E(X(t) | \mathcal{F}_s) = X(s) \text{ P a.e. } \omega.$$

Recalling the definition of conditional expectation, this says that for $F \in \mathcal{F}_s$

$$\int_F X(t) dP = \int_F E(X(t) | \mathcal{F}_s) dP = \int_F X(s) dP$$

for all $F \in \mathcal{F}_s$. A real valued stochastic process is called a sub-martingale if whenever $s \leq t$,

$$E(X(t) | \mathcal{F}_s) \geq X(s) \text{ a.e.}$$

and a supermartingale if

$$E(X(t) | \mathcal{F}_s) \leq X(s) \text{ a.e.}$$

Example 30.4.2 Let \mathcal{F}_t be a filtration and let Z be in $L^1(\Omega, \mathcal{F}_T, P)$. Then let $X(t) \equiv E(Z | \mathcal{F}_t)$.

This works because for $s < t$, $E(X(t) | \mathcal{F}_s) \equiv E(E(Z | \mathcal{F}_t) | \mathcal{F}_s) = E(Z | \mathcal{F}_s) \equiv X(s)$.

Proposition 30.4.3 The following statements hold for a stochastic process defined on $[0, T] \times \Omega$ having values in a real separable Banach space, E .

1. If $X(t)$ is a martingale then $\|X(t)\|, t \in [0, T]$ is a sub-martingale.

2. If g is an increasing convex function from $[0, \infty)$ to $[0, \infty)$ and

$$E(g(\|X(t)\|)) < \infty$$

for all $t \in [0, T]$ then $g(\|X(t)\|), t \in [0, T]$ is a sub-martingale.

Proof: Let $s \leq t$. Then from properties of conditional expectation and Theorem 24.12.1 on Page 702,

$$\begin{aligned} \|X(s)\| &= \|E(X(s) - X(t) | \mathcal{F}_s) + E(X(t) | \mathcal{F}_s)\| \\ &\stackrel{=0 \text{ a.e.}}{\leq} \overbrace{\|E(X(s) - X(t) | \mathcal{F}_s)\|}^{=0 \text{ a.e.}} + \|E(X(t) | \mathcal{F}_s)\| \leq \|E(X(t) | \mathcal{F}_s)\| \\ &\leq E(\|X(t)\| | \mathcal{F}_s) \end{aligned}$$

Consider the second claim. Recall Jensen's inequality for sub-martingales, Theorem 29.1.7 on Page 784. From the first part

$$\|X(s)\| \leq E(\|X(t)\| | \mathcal{F}_s) \text{ a.e.}$$

and so from Jensen's inequality,

$$g(\|X(s)\|) \leq g(E(\|X(t)\| | \mathcal{F}_s)) \leq E(g(\|X(t)\|) | \mathcal{F}_s) \text{ a.e.},$$

showing that $g(\|X(t)\|)$ is also a sub-martingale. ■

30.5 Some Maximal Estimates

Martingales and sub-martingales have some very interesting maximal estimates. I will present some of these here. The proofs are fairly general. For convenience, assume each \mathcal{F}_t contains the sets of measure zero from \mathcal{F} . This is so that it suffices to assume $t \rightarrow X(t)(\omega)$ is right continuous off some set of measure zero. If it were right continuous for each ω , then it wouldn't matter. Actually, in this book, I will mainly be interested in continuous processes. It is also possible to show that for real valued processes, one can get a right continuous version but this will not be used.

Lemma 30.5.1 *Let $\{\mathcal{F}_t\}$ be a filtration and let $\{X(t)\}$ be a nonnegative valued sub-martingale for $t \in [S, T]$. Then for $\lambda > 0$ and any $p \geq 1$, if, for each t , A_t is a \mathcal{F}_t measurable subset of $[X(t) > \lambda]$, then*

$$P(A_t) \leq \frac{1}{\lambda^p} \int_{A_t} X(T)^p dP.$$

Proof: From Jensen's inequality,

$$\begin{aligned} \lambda^p P(A_t) &\leq \int_{A_t} X(t)^p dP \leq \int_{A_t} E(X(T) | \mathcal{F}_t)^p dP \\ &\leq \int_{A_t} E(X(T)^p | \mathcal{F}_t) dP = \int_{A_t} X(T)^p dP \quad \blacksquare \end{aligned}$$

The following theorem is the main result.

Theorem 30.5.2 Let $\{\mathcal{F}_t\}$ be a filtration and let $\{X(t)\}$ be a nonnegative valued right continuous¹ sub-martingale for $t \in [S, T]$. Then for all $\lambda > 0$ and $p \geq 1$, for

$$X^* \equiv \sup_{t \in [S, T]} X(t),$$

$$P([X^* > \lambda]) \leq \frac{1}{\lambda^p} \int_{\Omega} \mathcal{X}_{[X^* > \lambda]} X(T)^p dP$$

In the case that $p > 1$, it is also true that

$$E((X^*)^p) \leq \left(\frac{p}{p-1}\right) E(X(T)^p)^{1/p} (E((X^*)^p))^{1/p'}$$

Also there are no measurability issues related to the above $\sup_{t \in [S, T]} X(t) \equiv X^*$. If $X(t) \in L^p(\Omega)$ for each t , then

$$E((X^*)^p)^{1/p} \leq \left(\frac{p}{p-1}\right) E(X(T)^p)^{1/p}$$

Thus X^* is also in $L^p(\Omega)$.

Proof: Let $S \leq t_0^m < t_1^m < \dots < t_{N_m}^m = T$ where $t_{j+1}^m - t_j^m = (T - S) 2^{-m}$. First consider $m = 1$.

$$A_{t_0^1} \equiv \{\omega \in \Omega : X(t_0^1)(\omega) > \lambda\}, A_{t_1^1} \equiv \{\omega \in \Omega : X(t_1^1)(\omega) > \lambda\} \setminus A_{t_0^1}$$

$$A_{t_2^1} \equiv \{\omega \in \Omega : X(t_2^1)(\omega) > \lambda\} \setminus (A_{t_0^1} \cup A_{t_1^1}).$$

Do this type of construction for $m = 2, 3, 4, \dots$ yielding disjoint sets, $\{A_{t_j^m}\}_{j=0}^{2^m}$ whose union equals

$$\cup_{t \in D_m} [X(t) > \lambda]$$

where $D_m = \{t_j^m\}_{j=0}^{2^m}$. Thus $D_m \subseteq D_{m+1}$. Then also, $D \equiv \cup_{m=1}^{\infty} D_m$ is dense and countable. From Lemma 30.5.1,

$$\begin{aligned} P(\cup_{t \in D_m} [X(t) > \lambda]) &= P\left(\left[\sup_{t \in D_m} X(t) > \lambda\right]\right) = \sum_{j=0}^{2^m} P(A_{t_j^m}) \\ &\leq \frac{1}{\lambda^p} \sum_{j=0}^{2^m} \int_{A_{t_j^m}} \mathcal{X}_{[\sup_{t \in D_m} X(t) > \lambda]} X(T)^p dP \end{aligned} \quad (30.22)$$

$$\leq \frac{1}{\lambda^p} \int_{\Omega} \mathcal{X}_{[\sup_{t \in D_m} X(t) > \lambda]} X(T)^p dP \leq \frac{1}{\lambda^p} \int_{\Omega} \mathcal{X}_{[\sup_{t \in D} X(t) > \lambda]} X(T)^p dP.$$

Let $m \rightarrow \infty$ in the above to obtain

$$P(\cup_{t \in D} [X(t) > \lambda]) = P\left(\left[\sup_{t \in D} X(t) > \lambda\right]\right) \leq \frac{1}{\lambda^p} \int_{\Omega} \mathcal{X}_{[\sup_{t \in D} X(t) > \lambda]} X(T)^p dP. \quad (30.23)$$

¹ $t \rightarrow X(t)(\omega)$ is continuous from the right for a.e. ω .

From now on, we begin using the assumption that for a.e. $\omega \in \Omega$, $t \rightarrow X(t)(\omega)$ is right continuous. Then with this assumption of right continuity, the following claim holds.

$$\sup_{t \in [S, T]} X(t) \equiv X^* = \sup_{t \in D} X(t)$$

which verifies that X^* is measurable. Then from 30.23,

$$\begin{aligned} P([X^* > \lambda]) &= P\left(\left[\sup_{t \in D} X(t) > \lambda\right]\right) \\ &\leq \frac{1}{\lambda^p} \int_{\Omega} \mathcal{X}_{[\sup_{t \in D} X(t) > \lambda]} X(T)^p dP = \frac{1}{\lambda^p} \int_{\Omega} \mathcal{X}_{[X^* > \lambda]} X(T)^p dP \end{aligned}$$

Now consider the other inequality. Using the distribution function technique and the above estimate obtained in the first part, and earlier facts about the distribution function,

$$E((X^*)^p) = \int_0^\infty p\alpha^{p-1} P([X^* > \alpha]) d\alpha$$

Then using Lemma 29.3.13 to justify interchange in order of integration,

$$\begin{aligned} &\leq \int_0^\infty p\alpha^{p-1} \frac{1}{\alpha} \int_{\Omega} \mathcal{X}_{[X^* > \alpha]} X(T) dP d\alpha = p \int_{\Omega} \int_0^{X^*} \alpha^{p-2} d\alpha X(T) dP \\ &= \frac{p}{p-1} \int_{\Omega} (X^*)^{p-1} X(T) dP \leq \frac{p}{p-1} \left(\int_{\Omega} (X^*)^p \right)^{1/p'} \left(\int_{\Omega} X(T)^p \right)^{1/p} \\ &= \frac{p}{p-1} E(X(T)^p)^{1/p} E((X^*)^p)^{1/p'}. \end{aligned} \quad (30.24)$$

Now assume $X(t) \in L^p(\Omega)$. Returning to 30.22, and letting X_n^* be $\sup_{t \in D_n} X(t)$, this says that

$$P([X_n^* > \lambda]) \leq \frac{1}{\lambda^p} \int_{\Omega} \mathcal{X}_{[X_n^* > \lambda]} X(T)^p dP$$

Then X_n^* achieves its maximum at one of finitely many values for t on a suitable subset of Ω . Thus it makes sense to write $\int_{\Omega} (X_n^*)^p dP$. Now repeat the argument. This yields

$$E((X_n^*)^p) \leq \frac{p}{p-1} E(X(T)^p)^{1/p} E((X_n^*)^p)^{1/p'}$$

Dividing by $E((X_n^*)^p)^{1/p'}$, one obtains

$$(E((X_n^*)^p))^{1/p} \leq \frac{p}{p-1} E(X(T)^p)^{1/p}$$

Now let $n \rightarrow \infty$ and use the monotone convergence theorem. ■

If you assumed $t \rightarrow X(t)$ is lower semi-continuous instead of right continuous, it appears the above argument would also work.

With Theorem 30.5.2, here is an important maximal estimate for martingales having values in E , a real separable Banach space. In the following, either $t \rightarrow X(t)(\omega)$ is right continuous for all ω or for a.e. ω each \mathcal{F}_t contains the sets of measure zero.

Theorem 30.5.3 *Let $X(t)$ for $t \in I = [0, T]$ be an E valued right continuous martingale with respect to a filtration \mathcal{F}_t . Then for $p \geq 1$,*

$$P\left(\left[\sup_{t \in I} \|X(t)\| > \lambda\right]\right) \leq \frac{1}{\lambda^p} E(\|X(T)\|^p). \quad (30.25)$$

If $p > 1$,

$$E\left(\left(\sup_{t \in [S, T]} \|X(t)\|\right)^p\right) \leq \left(\frac{p}{p-1}\right) E(\|X(T)\|^p)^{1/p} E\left(\left(\sup_{t \in [S, T]} \|X(t)\|\right)^p\right)^{1/p'} \quad (30.26)$$

If, in addition, each $X(t) \in L^p(\Omega)$ for each t , then

$$E\left(\left(\sup_{t \in [S, T]} \|X(t)\|\right)^p\right)^{1/p} \leq \left(\frac{p}{p-1}\right) E(\|X(T)\|^p)^{1/p} \quad (30.27)$$

Proof: By Proposition 30.4.3 $\|X(t)\|, t \in I$ is a sub-martingale and so from Theorem 30.5.2, it follows 30.25 and 30.26 hold. 30.27 also holds from Theorem 30.5.2. You just apply that theorem to the sub-martingale $Z(t) \equiv \|X(t)\|$ and let $Z^*(t) = \sup_{s \in [S, T]} \|X(s)\|$.
■

Chapter 31

Optional Sampling Theorems

As with discrete martingales, there is a notion of stopping time and optional sampling theorems. These are considered by approximating with discrete stopping times. It is like the case of the integral where one uses step functions or simple functions to approximate a given function.

31.1 Review of Discrete Stopping Times

First it is necessary to define the notion of a stopping time. The following definition was discussed earlier in the context of discrete processes.

Definition 31.1.1 Let (Ω, \mathcal{F}, P) be a probability space and let $\{\mathcal{F}_n\}_{n=1}^{\infty}$ be an increasing sequence of σ algebras each contained in \mathcal{F} , called a discrete filtration. A stopping time is a measurable function, τ which maps Ω to \mathbb{N} ,

$$\tau^{-1}(A) \in \mathcal{F} \text{ for all } A \in \mathcal{P}(\mathbb{N}),$$

such that for all $n \in \mathbb{N}$,

$$[\tau \leq n] \in \mathcal{F}_n.$$

Note this is equivalent to saying

$$[\tau = n] \in \mathcal{F}_n$$

because

$$[\tau = n] = [\tau \leq n] \setminus [\tau \leq n-1].$$

For τ a stopping time define \mathcal{F}_τ as follows.

$$\mathcal{F}_\tau \equiv \{A \in \mathcal{F} : A \cap [\tau \leq n] \in \mathcal{F}_n \text{ for all } n \in \mathbb{N}\}$$

These sets in \mathcal{F}_τ are referred to as “prior” to τ .

It is clear that \mathcal{F}_τ is a σ algebra.

The most important example of a stopping time is the first hitting time.

Example 31.1.2 The first hitting time of an adapted process $X(n)$ of a Borel set G is a stopping time. This is defined as

$$\tau \equiv \min \{k : X(k) \in G\}$$

To see this, note that

$$[\tau = n] = \cap_{k < n} [X(k) \in G^c] \cap [X(n) \in G] \in \mathcal{F}_n.$$

This led to the following proposition. It was Proposition 29.4.4.

Proposition 31.1.3 For τ a stopping time, \mathcal{F}_τ is a σ algebra and if $Y(k)$ is \mathcal{F}_k measurable for all k , $Y(k)$ having values in a separable Banach space E , then

$$\omega \rightarrow Y(\tau(\omega))$$

is \mathcal{F}_τ measurable.

To see this,

$$(Y \circ \tau)^{-1}(G) \cap [\tau \leq n] = \cup_k \left(\overbrace{Y(k)^{-1}(G)}^{\in \mathcal{F}_k} \right) \cap [\tau = k] \cap [\tau \leq n]$$

The term in the union is \emptyset if $k > n$ and so the whole thing reduces to

$$\cup_{k=1}^n \left(\overbrace{Y(k)^{-1}(G)}^{\in \mathcal{F}_k} \right) \cap [\tau = k] \in \mathcal{F}_n$$

showing that $(Y \circ \tau)^{-1}(G) \in \mathcal{F}_\tau$.

The following lemma contains the fundamental properties of stopping times for discrete filtrations. It was Lemma 29.4.7.

Lemma 31.1.4 *In the situation of Definition 31.1.1,*

1. $\mathcal{F}_\tau \cap [\tau = i] = \mathcal{F}_i \cap [\tau = i]$ and $E(X|\mathcal{F}_\tau) = E(X|\mathcal{F}_i)$ a.e. on the set $[\tau = i]$. Also if $A \in \mathcal{F}_\tau$ or \mathcal{F}_i , then $A \cap [\tau = i] \in \mathcal{F}_i \cap \mathcal{F}_\tau$.
2. $E(X|\mathcal{F}_\tau) = E(X|\mathcal{F}_i)$ a.e. on the set $[\tau \leq i]$.
3. Also, if $\sigma \leq \tau$, then $\mathcal{F}_\sigma \subseteq \mathcal{F}_\tau$

Proof: The first two are in the above mentioned lemma. The first part of 1. comes fairly quickly from the definition. The next part of 1. about the conditional expectations is essentially because one can regard $\mathcal{F}_\tau \cap [\tau = i]$ and $\mathcal{F}_i \cap [\tau = i]$ as two equal σ algebras contained in $[\tau = i]$ and so the two conditional expectations are the same on $[\tau = i]$. The third part of 1. also follows from the definition. Then 2. is clearly true from 1. applied to $[\tau = j]$ for $j \leq i$.

Say $A \in \mathcal{F}_\tau$. Then for $j \leq i$, $[\tau = j] \in \mathcal{F}_\tau$ because $[\tau = j] \cap [\tau \leq k] \in \mathcal{F}_k$ for each k . Thus

$$\begin{aligned} \int_{A \cap [\tau=j]} X dP &= \int_{A \cap [\tau=j]} E(X|\mathcal{F}_\tau) dP = \int_{\substack{A \cap [\tau=j] \\ \in \mathcal{F}_j}} E(X|\mathcal{F}_j) dP \\ &= \int_{A \cap [\tau=j]} E(X|\mathcal{F}_i) dP \end{aligned}$$

Since A is arbitrary, $E(X|\mathcal{F}_\tau) = E(X|\mathcal{F}_i)$.

Now consider 3. If $A \in \mathcal{F}_\sigma$, this means $A \cap [\sigma \leq i] \in \mathcal{F}_i$ or equivalently, $A \cap [\sigma = i] \in \mathcal{F}_i$ for all i . Take such an A . Then $A \cap [\tau = n] = \cup_{i=1}^n A \cap [\sigma = i] \in \mathcal{F}_n$ and so $\mathcal{F}_\sigma \subseteq \mathcal{F}_\tau$. ■

The assertion that

$$E(Y|\mathcal{F}_\tau) = E(Y|\mathcal{F}_k) \text{ a.e.}$$

on $[\tau = k]$ and that a function g which is \mathcal{F}_τ or \mathcal{F}_k measurable when restricted to $[\tau = k]$ is \mathcal{G} measurable for

$$\mathcal{G} = [\tau = k] \cap \mathcal{F}_\tau = [\tau = k] \cap \mathcal{F}_k$$

is the main result in the above lemma and this fact leads to the amazing Doob optional sampling theorem below. Also note that if $Y(k)$ is any process defined on the positive integers k , then by definition, $Y(k)(\omega) = Y(\tau(\omega))(\omega)$ on the set $[\tau = k]$ because τ is constant on this set.

31.2 Review of Doob Optional Sampling Theorem

With this lemma, here is a major theorem, the optional sampling theorem of Doob. This one is for martingales having values in a Banach space. To begin with, consider the case of a martingale defined on a countable set. This was discussed earlier but it is the sort of thing that seems to me should be repeated because it is so amazing.

Theorem 31.2.1 *Let $\{M(k)\}$ be a martingale having values in E a separable real Banach space with respect to the increasing sequence of σ algebras, $\{\mathcal{F}_k\}$ and let σ, τ be two stopping times such that τ is bounded. Then $M(\tau)$ defined as $\omega \rightarrow M(\tau(\omega))$ is integrable and*

$$M(\sigma \wedge \tau) = E(M(\tau) | \mathcal{F}_\sigma).$$

Proof: By Proposition 31.1.3 $M(\tau)$ is \mathcal{F}_τ measurable.

Next note that since τ is bounded by some l ,

$$\int_{\Omega} \|M(\tau(\omega))\| dP \leq \sum_{i=1}^l \int_{[\tau=i]} \|M(i)\| dP < \infty.$$

This proves the first assertion and makes possible the consideration of conditional expectation.

$(E(M(l) | \mathcal{F}_\tau) = M(\tau))$ Let $l \geq \tau$ as described above. Then for $k \leq l$, by Lemma 31.1.4,

$$\mathcal{F}_k \cap [\tau = k] = \mathcal{F}_\tau \cap [\tau = k] \equiv \mathcal{G}$$

implying that if g is either \mathcal{F}_k measurable or \mathcal{F}_τ measurable, then its restriction to $[\tau = k]$ is \mathcal{G} measurable and so if $A \in \mathcal{F}_\tau \cap [\tau = k]$ then

$$\begin{aligned} \int_A E(M(l) | \mathcal{F}_\tau) dP &\equiv \int_A M(l) dP = \int_A E(M(l) | \mathcal{F}_k) dP \\ &= \int_A M(k) dP = \int_A M(\tau) dP \text{ (on } A, \tau = k) \end{aligned}$$

Therefore, since A was arbitrary, $E(M(l) | \mathcal{F}_\tau) = M(\tau)$ a.e. on $[\tau = k]$ for every $k \leq l$. It follows $E(M(l) | \mathcal{F}_\tau) = M(\tau)$ a.e. since it is true on each $[\tau = k]$ for all $k \leq l$.

$(M(\sigma \wedge \tau) = E(M(\tau) | \mathcal{F}_\sigma))$ Now consider $E(M(\tau) | \mathcal{F}_\sigma)$ on the set $[\sigma = i] \cap [\tau = j]$. By Lemma 31.1.4, on this set,

$$E(M(\tau) | \mathcal{F}_\sigma) = E(M(\tau) | \mathcal{F}_i) = E(E(M(l) | \mathcal{F}_\tau) | \mathcal{F}_i) = E(E(M(l) | \mathcal{F}_j) | \mathcal{F}_i)$$

If $j \leq i$, this reduces to

$$E(M(l) | \mathcal{F}_j) = M(j) = M(\sigma \wedge \tau).$$

If $j > i$, this reduces to

$$E(M(l) | \mathcal{F}_i) = M(i) = M(\sigma \wedge \tau)$$

and since this exhausts all possibilities for values of σ and τ , it follows

$$E(M(\tau) | \mathcal{F}_\sigma) = M(\sigma \wedge \tau) \text{ a.e. } \blacksquare$$

You can also give a version of the above to sub-martingales. This requires the following very interesting decomposition of a sub-martingale into the sum of an increasing stochastic process and a martingale. This was presented earlier as Lemma 29.4.9.

Theorem 31.2.2 *Let $\{X_n\}$ be a sub-martingale. Then there exists a unique stochastic process, $\{A_n\}$ and martingale, $\{M_n\}$ such that*

1. $A_n(\omega) \leq A_{n+1}(\omega)$, $A_1(\omega) = 0$,
2. A_n is \mathcal{F}_{n-1} adapted for all $n \geq 1$ where $\mathcal{F}_0 \equiv \mathcal{F}_1$.

and also $X_n = M_n + A_n$.

Recall that the thing which works is

$$A(n) \equiv \sum_{k=0}^{n-1} E(X(k+1) - X(k) | \mathcal{F}_k), \quad A(0) = 0$$

and that this is the only thing which will do what is required.

Now here is a version of the optional sampling theorem for sub-martingales. This was also presented earlier. However, it is good to go through the proof as a review.

Theorem 31.2.3 *Let $\{X(k)\}$ be a real valued sub-martingale with respect to the increasing sequence of σ algebras, $\{\mathcal{F}_k\}$ such that τ is bounded. Then $X(\tau)$ defined as*

$$\omega \rightarrow X(\tau(\omega))$$

is integrable and

$$X(\sigma \wedge \tau) \leq E(X(\tau) | \mathcal{F}_\sigma)$$

Proof: That $\omega \rightarrow X(\tau(\omega))$ is integrable follows right away as in the optional sampling theorem for martingales. You just consider the finitely many values of τ .

Use Theorem 31.2.2 above to write

$$X(n) = M(n) + A(n)$$

where M is a martingale and A is increasing with $A(n)$ being \mathcal{F}_{n-1} measurable and $A(0) = 0$ as discussed in Theorem 31.2.2. Then

$$E(X(\tau) | \mathcal{F}_\sigma) = E(M(\tau) + A(\tau) | \mathcal{F}_\sigma)$$

Now since A is increasing, you can use the optional sampling theorem for martingales to conclude that, since $\mathcal{F}_{\sigma \wedge \tau} \subseteq \mathcal{F}_\sigma$ and $A(\sigma \wedge \tau)$ is $\mathcal{F}_{\sigma \wedge \tau}$ measurable,

$$\begin{aligned} &\geq E(M(\tau) + A(\sigma \wedge \tau) | \mathcal{F}_\sigma) = E(M(\tau) | \mathcal{F}_\sigma) + A(\sigma \wedge \tau) \\ &= M(\sigma \wedge \tau) + A(\sigma \wedge \tau) = X(\sigma \wedge \tau). \blacksquare \end{aligned}$$

Note that if $\sigma \leq \tau$, the conclusion is $X(\sigma) \leq E(X(\tau) | \mathcal{F}_\sigma)$.

31.3 Doob Optional Sampling Continuous Case

31.3.1 Stopping Times

Let $X(t)$ be a stochastic process adapted to a filtration $\{\mathcal{F}_t\}$ for $t \in [0, T]$ meaning that $X(t)$ is \mathcal{F}_t measurable each \mathcal{F}_t being a subset of \mathcal{F} . We will assume two things. The stochastic process is right continuous and the filtration is normal. Recall what this means:

Definition 31.3.1 A normal filtration is one which satisfies the following :

1. \mathcal{F}_0 contains all $A \in \mathcal{F}$ such that $P(A) = 0$. Here \mathcal{F} is the σ algebra which contains all \mathcal{F}_t .
2. $\mathcal{F}_t = \mathcal{F}_{t+}$ for all $t \in I$ where $\mathcal{F}_{t+} \equiv \bigcap_{s>t} \mathcal{F}_s$.

For an \mathcal{F} measurable $[0, \infty)$ valued function τ to be a stopping time, we want to have the stopped process X^τ defined by $X^\tau(t)(\omega) \equiv X(t \wedge \tau(\omega))(\omega)$ to be adapted whenever X is right continuous and adapted. Thus a stopping time is a measurable function which can be used to stop the process while retaining the property of being adapted. The definition of such a condition which will make τ a stopping time is the same as in the case of a discreet process.

Definition 31.3.2 τ an \mathcal{F} measurable function is a stopping time if $[\tau \leq t] \in \mathcal{F}_t$.

Then this definition does what is desired. This is in the following proposition. For convenience, here is a definition.

Definition 31.3.3 Let $\{t\}_k \equiv 2^{-k}n$ where n is as large as possible and have $2^{-k}n \leq t$.

It seems like the theory is based on reducing to discrete stopping times defined as follows.

Definition 31.3.4

$$\tau_k(\omega) \equiv \sum_{n=0}^{\infty} \mathcal{X}_{\tau^{-1}((n2^{-k}, (n+1)2^{-k}])}(\omega) (n+1)2^{-k}.$$

Thus τ_k has values in the set $\{n2^{-k}\}_{n=0}^{\infty}$, $\tau_k \geq \tau$ and τ_k is within 2^{-k} of τ .

Then τ_k is a discrete stopping time with respect to the increasing σ algebras $\mathcal{F}_{\{t\}_k}$. This is in the following lemma.

Lemma 31.3.5 Let τ be a stopping time and let τ_k be defined above in Definition 31.3.4. Then $[\tau_k \leq \{t\}_k] \in \mathcal{F}_{\{t\}_k}$ and for all t , $[\tau_k \leq t] \in \mathcal{F}_t$. If you have finitely many stopping times $\{\sigma^k\}_{k=1}^n$ for $n < \infty$ then $\sigma \equiv \min\{\sigma^k\}_{k=1}^n$ and $\sigma \equiv \max\{\sigma^k\}_{k=1}^n$ are also stopping times.

Proof: If $t = (n+1)2^{-k}$ then $[\tau_k \leq \{t\}_k]$ is the same as $[\tau \leq (n+1)2^{-k}] = [\tau \leq t] \in \mathcal{F}_t = \mathcal{F}_{\{t\}_k}$. The other case is where for some n , $t \in (n2^{-k}, (n+1)2^{-k})$. Then $\{t\}_k = n2^{-k}$ and so $[\tau_k \leq \{t\}_k]$ is $[\tau \leq n2^{-k}] \in \mathcal{F}_{n2^{-k}} = \mathcal{F}_{\{t\}_k}$. Thus, in particular, $[\tau_k \leq t] \in \mathcal{F}_t$ since $\{t\}_k \leq t$.

As to the second claim,

$$\left[\min\{\sigma^k\}_{k=1}^n \leq t \right] = \bigcup_{k=1}^n [\sigma^k \leq t] \in \mathcal{F}_t$$

and $[\max\{\sigma^k\}_{k=1}^n \leq t] = \bigcap_{k=1}^n [\sigma^k \leq t] \in \mathcal{F}_t$. ■

Proposition 31.3.6 *Let $\{\mathcal{F}_t\}$ be a normal filtration and let $X(t)$ be a right continuous process adapted to $\{\mathcal{F}_t\}$. Then if τ is a stopping time, it follows that the stopped process X^τ defined by $X^\tau(t) \equiv X(\tau \wedge t)$ is also adapted.*

Proof: Let $\tau_k(\omega) \equiv \sum_{n=0}^{\infty} \mathcal{X}_{\tau^{-1}((n2^{-k}, (n+1)2^{-k}])}(\omega)(n+1)2^{-k}$. Thus τ_k has discrete values $n2^{-k}, n = 1, \dots$ and $\tau_k(\omega) = (n+1)2^{-k}$ exactly when ω is in

$$\left[\tau \in [0, (n+1)2^{-k}] \right] \setminus \left[\tau \in [0, n2^{-k}] \right]$$

and these sets are both in $\mathcal{F}_{(n+1)2^{-k}}$. Now consider $X(\tau \wedge t)^{-1}(O)$ for O an open set. Since O is open, it follows from right continuity of X that if $X(\tau \wedge t) \in O$, then $X(\tau_k \wedge t) \in O$ whenever k is large enough, depending on ω of course. Thus

$$X(\tau \wedge t)^{-1}(O) = \bigcup_{m=1}^{\infty} \bigcap_{k \geq m} X(\tau_k \wedge t)^{-1}(O)$$

Now $X(\tau_k \wedge t)^{-1}(O) = \left([\tau_k \leq \{t\}_k] \cap X(\tau_k)^{-1}(O) \right) \cup \left([\tau_k > \{t\}_k] \cap X(t)^{-1}(O) \right)$. The second term in the union is in \mathcal{F}_t because X is adapted and $[\tau_k > \{t\}_k]$ is the complement of the set $[\tau_k \leq \{t\}_k]$ which is in \mathcal{F}_t . The first term is of the form $[\tau_k \leq \{t\}_k] \cap \left(\bigcup_{j=0}^{2^k \{t\}_k} X(2^{-k}j)^{-1}(O) \right) \in \mathcal{F}_t$ again because X is adapted. Since each $X(\tau_k \wedge t)^{-1}(O) \in \mathcal{F}_t$, it follows that $X(\tau \wedge t)^{-1}(O) = \bigcup_{m=1}^{\infty} \bigcap_{k \geq m} X(\tau_k \wedge t)^{-1}(O)$ is in \mathcal{F}_t . ■

By analogy to the discrete case, here are the prior sets.

Definition 31.3.7 *Let (Ω, \mathcal{F}, P) be a probability space and let \mathcal{F}_t be a filtration. Recall a measurable function, $\tau : \Omega \rightarrow [0, \infty]$ is called a stopping time if*

$$[\tau \leq t] \in \mathcal{F}_t$$

for all $t \geq 0$. Associated with a stopping time is the σ algebra, \mathcal{F}_τ defined by

$$\mathcal{F}_\tau \equiv \{A \in \mathcal{F} : A \cap [\tau \leq t] \in \mathcal{F}_t \text{ for all } t\}.$$

These sets are also called those “prior” to τ .

Note that \mathcal{F}_τ is obviously closed with respect to countable unions. If $A \in \mathcal{F}_\tau$, then

$$A^C \cap [\tau \leq t] = [\tau \leq t] \setminus (A \cap [\tau \leq t]) \in \mathcal{F}_t$$

Thus \mathcal{F}_τ is a σ algebra. What if $\tau \leq \sigma$? Does it follow that $\mathcal{F}_\tau \subseteq \mathcal{F}_\sigma$? What about $\omega \rightarrow X(\tau)$? Is this measurable?

Recall $\{t\}_k \equiv n2^{-k}$ where n is as large as possible with $n2^{-k} \leq t$.

Proposition 31.3.8 *Let τ be a stopping time and let τ_k be the stopping time having discrete values described in Lemma 31.3.5 which for a given stopping time τ is given by*

$$\tau_k(\omega) \equiv \sum_{n=0}^{\infty} \mathcal{X}_{\tau^{-1}((n2^{-k}, (n+1)2^{-k}])}(\omega)(n+1)2^{-k}.$$

Then

1. $\mathcal{F}_\tau \subseteq \mathcal{F}_{\tau_k}$ and if \mathcal{F}_t is normal, then if $A \in \mathcal{F}_{\tau_k}$ for all k , it follows that $A \in \mathcal{F}_\tau$.
2. $A \in \mathcal{F}_{\tau_k}$ if and only if $A \cap [\tau_k = n2^{-k}] \in \mathcal{F}_{n2^{-k}}$ for all n .
3. If $\tau \leq \sigma$, then $\tau_k \leq \sigma_k$.
4. More generally, if $\tau \leq \sigma$ are two stopping times, then $\mathcal{F}_\tau \subseteq \mathcal{F}_\sigma$.
5. If $X(t)$ is a right continuous process adapted to the normal filtration \mathcal{F}_t and τ is a stopping time, then $\omega \rightarrow X(\tau(\omega))$ is \mathcal{F}_τ measurable. Here X has values in a Banach space.

Proof: 1.) Let $A \in \mathcal{F}_\tau$. Then for $t \in (n2^{-k}, (n+1)2^{-k}]$, if $t < (n+1)2^{-k}$,

$$A \cap [\tau_k \leq t] = A \cap [\tau_k \leq n2^{-k}] = A \cap [\tau \leq n2^{-k}] \in \mathcal{F}_{n2^{-k}} \subseteq \mathcal{F}_t$$

so $A \in \mathcal{F}_{\tau_k}$. If $t = (n+1)2^{-k}$ then $A \cap [\tau_k \leq t] \in \mathcal{F}_{(n+1)2^{-k}} = \mathcal{F}_t$. Thus $A \in \mathcal{F}_{\tau_k}$ and $\mathcal{F}_\tau \subseteq \mathcal{F}_{\tau_k}$. Now for the other part, Consider $A \cap [\tau \leq t]$. If $t \in (n2^{-k}, (n+1)2^{-k}]$, then $[\tau \leq t] = [\tau_k \leq (n+1)2^{-k}]$ and so $A \cap [\tau \leq t] = A \cap [\tau_k \leq (n+1)2^{-k}] \in \mathcal{F}_{(n+1)2^{-k}} \subseteq \mathcal{F}_{t+2^{-k}}$. Since the filtration is normal, and for all k , $A \cap [\tau \leq t] \in \mathcal{F}_{t+2^{-k}}$, it follows that $A \cap [\tau \leq t] \in \mathcal{F}_t$ and so $A \in \mathcal{F}_\tau$ as claimed.

2.) Note

$$\begin{aligned} [\tau_k = n2^{-k}] &= [\tau_k \leq n2^{-k}] \setminus [\tau_k \leq (n-1)2^{-k}] \\ &= [\tau \leq 2^{-k}n] \setminus [\tau \leq (n-1)2^{-k}] \end{aligned}$$

If $A \in \mathcal{F}_{\tau_k}$ then by definition,

$$\begin{aligned} A \cap [\tau_k = n2^{-k}] &= A \cap ([\tau_k \leq n2^{-k}] \setminus [\tau_k \leq (n-1)2^{-k}]) \\ &= A \cap ([\tau \leq 2^{-k}n] \setminus [\tau \leq (n-1)2^{-k}]) \in \mathcal{F}_{n2^{-k}} \end{aligned}$$

Conversely, if

$$A \cap [\tau_k = n2^{-k}] \in \mathcal{F}_{n2^{-k}}$$

for all n , then consider $A \cap [\tau_k \leq t]$ for $t \in (n2^{-k}, (n+1)2^{-k}]$. Is $A \cap [\tau_k \leq t] \in \mathcal{F}_t$? If $t < (n+1)2^{-k}$, then $A \cap [\tau_k \leq t] = \sum_{j=1}^n A \cap [\tau_k = j2^{-k}] \in \mathcal{F}_{n2^{-k}} \subseteq \mathcal{F}_t$. If $t = (n+1)2^{-k}$, $A \cap [\tau_k \leq t] = \sum_{j=1}^{n+1} A \cap [\tau_k = j2^{-k}] \in \mathcal{F}_t$ and so $A \in \mathcal{F}_{(n+1)2^{-k}} \subseteq \mathcal{F}_{\tau_k}$. This proves 2.).

3.) The values of both τ_k and σ_k are $n2^{-k}$ for some nonnegative integer n . τ_k equals $n2^{-k}$ on $\tau^{-1}((n-1)2^{-k}, n2^{-k})$. Thus on this set, σ cannot be smaller than or equal to $(n-1)2^{-k}$. Hence σ_k is at least $n2^{-k}$.

4.) Let $\tau \leq \sigma$ and let $A \in \mathcal{F}_\tau$. $A \cap [\sigma \leq t] = A \cap [\tau \leq t] \cap [\sigma \leq t]$ because $\tau \leq \sigma$. However, $A \cap [\tau \leq t] \in \mathcal{F}_t$ and $[\sigma \leq t] \in \mathcal{F}_t$ so the right side is in \mathcal{F}_t which means $A \in \mathcal{F}_\sigma$.

5.) Let U be an open set.

$$X(\tau_k)^{-1}(U) \cap [\tau_k < t] = \cup_{j2^{-k} \leq \{t\}_k} X(j2^{-k})^{-1}(U) \in \mathcal{F}_{\{t\}_k}.$$

Now say $t \in ((n-1)2^{-k}, n2^{-k})$. If $t < (n+1)2^{-k}$, then $X(\tau_k)^{-1}(U) \cap [\tau_k \leq t]$ was just given. It is in $\mathcal{F}_{n2^{-k}} = \mathcal{F}_{\{t\}_k} \subseteq \mathcal{F}_t$. Otherwise $t = (n+1)2^{-k}$ and in this case,

$X(\tau_k)^{-1}(U) \cap [\tau_k \leq t]$ is in $\mathcal{F}_{(n+1)2^{-k}} = \mathcal{F}_t$. Thus $\omega \rightarrow X(\tau_k)(\omega)$ is \mathcal{F}_{τ_k} measurable. It follows that for a fixed \hat{k} , $X(\tau_k)$ is $\mathcal{F}_{\tau_{\hat{k}}}$ measurable for each $k > \hat{k}$ because τ_k is decreasing in k so this follows from Part 4. Now let $k \rightarrow \infty$ and use right continuity of X to conclude that $X(\tau)$ is $\mathcal{F}_{\tau_{\hat{k}}}$ measurable. Thus $X(\tau)^{-1}(U) \in \mathcal{F}_{\tau_{\hat{k}}}$ for each \hat{k} and so, by Part 1, it follows that $X(\tau)^{-1}(U) \in \mathcal{F}_{\tau}$. Therefore, $X(\tau)$ is \mathcal{F}_{τ} measurable. ■

Next is an important proposition which gives a typical example of a stopping time. Since the process has t in an interval, one must be much more careful about the nature of the set which is hit.

Proposition 31.3.9 *Let B be an open subset of topological space E and let $X(t)$ be a right continuous \mathcal{F}_t adapted stochastic process such that \mathcal{F}_t is normal. Then define*

$$\tau(\omega) \equiv \inf\{t > 0 : X(t)(\omega) \in B\}.$$

This is called the first hitting time. Then τ is a stopping time. If $X(t)$ is continuous and adapted to \mathcal{F}_t , a normal filtration, then if H is a nonempty closed set such that $H = \bigcap_{n=1}^{\infty} B_n$ for B_n open, $B_n \supseteq B_{n+1}$,

$$\tau(\omega) \equiv \inf\{t > 0 : X(t)(\omega) \in H\}$$

is also a stopping time.

Proof: Consider the first claim. $\omega \in [\tau = a]$ implies that for each $n \in \mathbb{N}$, there exists $t \in [a, a + \frac{1}{n}]$ such that $X(t) \in B$. Also for $t < a$, you would need $X(t) \notin B$. By right continuity, this is the same as saying that $X(d) \notin B$ for all rational $d < a$. (If $t < a$, you could let $d_n \downarrow t$ where $X(d_n) \in B^C$, a closed set. Then it follows that $X(t)$ is also in the closed set B^C .) Thus, aside from a set of measure zero, for each $m \in \mathbb{N}$,

$$[\tau = a] = \left(\bigcap_{n=m}^{\infty} \bigcup_{t \in [a, a + \frac{1}{n}]} [X(t) \in B] \right) \cap \left(\bigcap_{t \in [0, a)} [X(t) \in B^C] \right)$$

Since $X(t)$ is right continuous, this is the same as

$$\left(\bigcap_{n=m}^{\infty} \bigcup_{d \in \mathbb{Q} \cap [a, a + \frac{1}{n}]} [X(d) \in B] \right) \cap \left(\bigcap_{d \in \mathbb{Q} \cap [0, a)} [X(d) \in B^C] \right) \in \mathcal{F}_{a + \frac{1}{m}}$$

Thus, since the filtration is normal,

$$[\tau = a] \in \bigcap_{m=1}^{\infty} \mathcal{F}_{a + \frac{1}{m}} = \mathcal{F}_{a+} = \mathcal{F}_a$$

I want to consider $[\tau \leq a]$. What of $[\tau < a]$? This is equivalent to saying that $X(t) \in B$ for some $t < a$. Since X is right continuous, this is the same as saying that $X(t) \in B$ for some $t \in \mathbb{Q}, t < a$. Thus

$$[\tau < a] = \bigcup_{d \in \mathbb{Q}, d < a} [X(d) \in B] \in \mathcal{F}_a$$

It follows that $[\tau \leq a] = [\tau < a] \cup [\tau = a] \in \mathcal{F}_a$. Thus τ is indeed a stopping time.

Now consider the claim involving the additional assumption that $X(t)$ is continuous and it is desired to hit a closed set $H = \bigcap_{n=1}^{\infty} B_n$ where B_n is open, $B_n \supseteq B_{n+1}$. (Note that if the topological space is a metric space, this is always possible so this is not a big restriction.) Then let τ_n be the first hitting time of B_n by $X(t)$. Then it can be shown that

$$[\tau \leq a] = \bigcap_n [\tau_n \leq a] \in \mathcal{F}_a$$

To show this, first note that $\omega \in [\tau \leq a]$ if and only if there exists $t \leq a$ such that $X(t)(\omega) \in H$. This follows from continuity and the fact that H is closed. Thus $\tau_n \leq a$ for all n because for some $t \leq a$, $X(t) \in H \subseteq B_n$ for all n . Next suppose $\omega \in [\tau_n \leq a]$ for all n . Then for $\delta_n \downarrow 0$, there exists $t_n \in [0, a + \delta_n]$ such that $X(t_n)(\omega) \in B_n$. It follows there is a subsequence, still denoted by t_n such that $t_n \rightarrow t \in [0, a]$. By continuity of X , it must be the case that $X(t)(\omega) \in H$ and so $\omega \in [\tau \leq a]$. This shows the above formula. Now by the first part, each $[\tau_n \leq a] \in \mathcal{F}_a$ and so $[\tau \leq a] \in \mathcal{F}_a$ also. ■

Another useful result for real valued stochastic process is the following in which continuity is generalized to lower semicontinuity.

Proposition 31.3.10 *Let $X(t)$ be a real valued stochastic process which is \mathcal{F}_t adapted for a normal filtration \mathcal{F}_t , with the property that $t \rightarrow X(t)$ is lower semicontinuous. Then*

$$\tau \equiv \inf\{t : X(t) > \alpha\}$$

is a stopping time.

Proof: As above, for each $m > 0$,

$$[\tau = a] = \left(\bigcap_{n=m}^{\infty} \bigcup_{t \in [a, a + \frac{1}{n}]} [X(t) > \alpha] \right) \cap \left(\bigcap_{t \in [0, a)} [X(t) \leq \alpha] \right)$$

Now

$$\bigcap_{t \in [0, a)} [X(t) \leq \alpha] \subseteq \bigcap_{t \in [0, a), t \in \mathbb{Q}} [X(t) \leq \alpha]$$

If ω is in the right side, then for arbitrary $t < a$, let $t_n \downarrow t$ where $t_n \in \mathbb{Q}$ and $t < a$. Then $X(t) \leq \liminf_{n \rightarrow \infty} X(t_n) \leq \alpha$ and so ω is in the left side also. Thus

$$\bigcap_{t \in [0, a)} [X(t) \leq \alpha] = \bigcap_{t \in [0, a), t \in \mathbb{Q}} [X(t) \leq \alpha]$$

$$\bigcup_{t \in [a, a + \frac{1}{n}]} [X(t) > \alpha] \supseteq \bigcup_{t \in [a, a + \frac{1}{n}], t \in \mathbb{Q}} [X(t) > \alpha]$$

If ω is in the left side, then for some t in the given interval, $X(t) > \alpha$. If for all $s \in [a, a + \frac{1}{n}] \cap \mathbb{Q}$ you have $X(s) \leq \alpha$, then you could take $s_n \rightarrow t$ where $X(s_n) \leq \alpha$ and conclude from lower semicontinuity that $X(t) \leq \alpha$ also. Thus there is some rational s where $X(s) > \alpha$ and so the two sides are equal. Hence,

$$[\tau = a] = \left(\bigcap_{n=m}^{\infty} \bigcup_{t \in [a, a + \frac{1}{n}], t \in \mathbb{Q}} [X(t) > \alpha] \right) \cap \left(\bigcap_{t \in [0, a), t \in \mathbb{Q}} [X(t) \leq \alpha] \right)$$

The first set on the right is in $\mathcal{F}_{a+(1/m)}$ and so is the next set on the right. Hence $[\tau = a] \in \bigcap_m \mathcal{F}_{a+(1/m)} = \mathcal{F}_a$. To be a stopping time, one needs $[\tau \leq a] \in \mathcal{F}_a$. What of $[\tau < a]$? This equals $\bigcup_{t \in [0, a)} [X(t) > \alpha] = \bigcup_{t \in [0, a) \cap \mathbb{Q}} [X(t) > \alpha] \in \mathcal{F}_a$, the equality following from lower semi-continuity. Thus $[\tau \leq a] = [\tau = a] \cup [\tau < a] \in \mathcal{F}_a$. ■

31.3.2 The Optional Sampling Theorem Continuous Case

Proposition 31.3.11 *Let $M(t), t \geq 0$ be a martingale with values in E a separable Banach space and let τ be a bounded stopping time whose maximum is T . Then $M(\tau)$ is \mathcal{F}_τ measurable and in fact, $\int_\Omega \|M(\tau)\| dP < \infty$. Letting*

$$\tau_k(\omega) \equiv \sum_{n=0}^{\infty} \mathcal{X}_{\tau^{-1}((n2^{-k}, (n+1)2^{-k}])}(\omega) (n+1) 2^{-k}$$

be the discrete stopping times. $\{M(\tau_k)\}$ are uniformly integrable on \mathcal{F}_T .

Proof: From Proposition 31.3.8 $M(\tau)$ is measurable. Since τ is bounded, this is always a finite sum for τ_k . Then each ω is in exactly one $\tau^{-1}((n2^{-k}, (n+1)2^{-k}])$ for some n . Say $\omega \in \tau^{-1}((n2^{-k}, (n+1)2^{-k}])$. Then for that ω , $M(\tau_k(\omega)) = M((n+1)2^{-k})(\omega)$. Thus $M(\tau_k(\omega))$ is given by

$$M(\tau_k(\omega)) = \sum_{n=0}^{\infty} \mathcal{X}_{\tau^{-1}((n2^{-k}, (n+1)2^{-k}])}(\omega) M((n+1)2^{-k})(\omega)$$

and

$$\begin{aligned} \|M(\tau_k)(\omega)\| &\leq \sum_{n=0}^{\infty} \mathcal{X}_{\tau^{-1}(I_n)}(\omega) \|M((n+1)2^{-k})\|, \\ I_n &\equiv (n2^{-k}, (n+1)2^{-k}] \end{aligned}$$

Then, since $M(t)$ is a martingale, $E(M((n+1)2^{-k}) | \mathcal{F}_{n2^{-k}}) = M(n2^{-k})$ and so

$$\begin{aligned} \|M(n2^{-k})\| &= \|E(M((n+1)2^{-k}) | \mathcal{F}_{n2^{-k}})\| \\ &\leq E(\|M((n+1)2^{-k})\| | \mathcal{F}_{n2^{-k}}). \end{aligned}$$

Therefore, iterating this gives

$$\|M((n+1)2^{-k})\| \leq E(\|M(T_k)\| | \mathcal{F}_{(n+1)2^{-k}})$$

where T_k is the smallest number greater than or equal to T which is of the form $m2^{-k}$ for m a positive integer. Then, since $\mathcal{X}_{\tau^{-1}(I_n)}$ is $\mathcal{F}_{(n+1)2^{-k}}$ measurable, it is \mathcal{F}_{T_k} measurable and so

$$\begin{aligned} \|M(\tau_k)\| &\leq \sum_{n=0}^{\infty} \mathcal{X}_{\tau^{-1}(I_n)}(\omega) E(\|M(T_k)\| | \mathcal{F}_{(n+1)2^{-k}}) \\ &= \sum_{n=0}^{\infty} E(\mathcal{X}_{\tau^{-1}(I_n)} \|M(T_k)\| | \mathcal{F}_{(n+1)2^{-k}}) \end{aligned}$$

because $\mathcal{X}_{\tau^{-1}(I_n)}$ is $\mathcal{F}_{(n+1)2^{-k}}$ measurable. Thus

$$\int_{\Omega} \|M(\tau_k)\| dP \leq \sum_{n=0}^{\infty} \int_{\Omega} \mathcal{X}_{\tau^{-1}(I_n)} \|M(T_k)\| dP = \int_{\Omega} \|M(T_k)\| dP.$$

Thus

$$\int_{\Omega} \|M(\tau_k)\| dP \leq \int_{\Omega} \|M(T_k)\| dP$$

Now use right continuity and Fatou's lemma.

$$\int_{\Omega} \|M(\tau)\| dP \leq \liminf_{k \rightarrow \infty} \int_{\Omega} \|M(\tau_k)\| dP \leq \liminf_{k \rightarrow \infty} \int_{\Omega} \|M(T_k)\| dP$$

Pick $\hat{T} > T_k$ for all $k = 1, 2, \dots$. Then

$$M(T_k) = E(M(\hat{T}) | \mathcal{F}_{T_k})$$

and so $\|M(T_k)\| \leq E(\|M(\hat{T})\| | \mathcal{F}_{T_k})$ and so

$$\int_{\Omega} \|M(T_k)\| dP \leq \int_{\Omega} E(\|M(\hat{T})\| | \mathcal{F}_{T_k}) dP = \int_{\Omega} \|M(\hat{T})\| dP < \infty$$

Therefore,

$$\int_{\Omega} \|M(\tau)\| dP \leq \liminf_{k \rightarrow \infty} \int_{\Omega} \|M(\tau_k)\| dP \leq \int_{\Omega} \|M(\hat{T})\| dP < \infty$$

because it is given that $M(t)$ is in L^1 for each t .

In the above, you could replace Ω with $A \in \mathcal{F}_T$ and conclude

$$\int_A \|M(\tau_k)\| dP \leq \int_A \|M(\hat{T})\| dP$$

which implies the $M(\tau_k)$ are uniformly integrable. Given $\varepsilon > 0$ there is $\delta > 0$ such that if $P(A) < \delta$ then

$$\int_A \|M(\hat{T})\| dP < \varepsilon$$

and so also $\int_A \|M(\tau_k)\| dP < \varepsilon$. ■

Now consider an increasing in t family of stopping times, $\tau(t) (\omega \rightarrow \tau(t)(\omega))$. It turns out this is a sub-martingale.

Lemma 31.3.12 *Let $\{\tau(t)\}$ be an increasing in t family of stopping times, $\tau(t) \geq \tau(s)$ if $s < t$. Then $\tau(t)$ is adapted to the σ algebras $\mathcal{F}_{\tau(t)}$ and $\{\tau(t)\}$ is a sub-martingale adapted to these σ algebras.*

Proof: First I need to show that a stopping time, τ is \mathcal{F}_{τ} measurable. Consider $[\tau \leq s]$. Is this in \mathcal{F}_{τ} ? Is $[\tau \leq s] \cap [\tau \leq r] \in \mathcal{F}_r$ for each r ? This is obviously so if $s \leq r$ because the intersection reduces to $[\tau \leq s] \in \mathcal{F}_s \subseteq \mathcal{F}_r$. On the other hand, if $s > r$ then the intersection reduces to $[\tau \leq r] \in \mathcal{F}_r$ and so it is clear that τ is \mathcal{F}_{τ} measurable. It remains to verify that $t \rightarrow \tau(t)$ is a sub-martingale.

Let $s < t$ and let $A \in \mathcal{F}_{\tau(s)}$

$$\int_A E(\tau(t) | \mathcal{F}_{\tau(s)}) dP \equiv \int_A \tau(t) dP \geq \int_A \tau(s) dP$$

and this shows $E(\tau(t) | \mathcal{F}_{\tau(s)}) \geq \tau(s)$ so this is a submartingale as claimed. ■

Now here is an important example. Recall that for τ a stopping time, so is $t \vee \tau$ because

$$[t \vee \tau \leq s] = \begin{cases} \Omega & \text{if } t \leq s, \\ \emptyset & \text{otherwise} \end{cases} \cap [\tau \leq s] \in \mathcal{F}_s.$$

Also recall that if σ is a stopping time, then for adapted Y , $Y(\sigma)$ is \mathcal{F}_{σ} adapted. This is in Proposition 31.3.8.

Proposition 31.3.13 *Let τ be a stopping time and let X be continuous and adapted to the filtration \mathcal{F}_t . Then for $a > 0$, define σ as*

$$\sigma(\omega) \equiv \inf\{t > \tau(\omega) : \|X(t)(\omega) - X(\tau(\omega))\| = a\}$$

Then σ is also a stopping time. That is $[\sigma \leq t] \in \mathcal{F}_t$.

Proof: To see this is so, let

$$Y(t)(\omega) = \|X(t \vee \tau)(\omega) - X(\tau(\omega))\|$$

Then $Y(t)$ is $\mathcal{F}_{t \vee \tau}$ measurable. It is desired to show that Y is \mathcal{F}_t adapted. Hence if U is open in \mathbb{R} , then

$$Y(t)^{-1}(U) = \left(Y(t)^{-1}(U) \cap [\tau \leq t] \right) \cup \left(Y(t)^{-1}(U) \cap [\tau > t] \right)$$

The second set in the above union on the right equals either \emptyset or $[\tau > t]$ depending on whether $0 \in U$. If $\tau > t$, then $Y(t) = 0$ and so the second set equals $[\tau > t]$ if $0 \in U$. If $0 \notin U$, then the second set equals \emptyset . Thus the second set above is in \mathcal{F}_t . It is necessary to show the first set is also in \mathcal{F}_t . The first set equals

$$Y(t)^{-1}(U) \cap [\tau \leq t] = Y(t)^{-1}(U) \cap [\tau \vee t \leq t]$$

because $[\tau \vee t \leq t] = [\tau \leq t]$. However, $Y(t)^{-1}(U) \in \mathcal{F}_{t \vee \tau}$ and so the set on the right in the above is in \mathcal{F}_t . For A to be in $\mathcal{F}_{t \vee \tau}$ means $A \cap [\tau \vee t \leq s] \in \mathcal{F}_s$ for each s . In particular, this is true for $s = t$. Therefore, $Y(t)$ is adapted. Then σ is just the first hitting time for $Y(t)$ to equal the closed set a . Therefore, σ is a stopping time by Proposition 31.3.9. ■

The following corollary involves the same argument. Just replace

$$\|X(t \vee \tau)(\omega) - X(\tau(\omega))\|$$

with $g(X(t \vee \tau)(\omega) - X(\tau(\omega)))$.

Corollary 31.3.14 *Let τ be a stopping time and let X be continuous and adapted to the filtration \mathcal{F}_t . Also let g be a continuous real valued function. Then for $a > 0$, define σ as*

$$\sigma(\omega) \equiv \inf \{t > \tau(\omega) : g(X(t)(\omega) - X(\tau(\omega))) = a\}$$

Then σ is also a stopping time.

Next I want a version of the Doob optional sampling theorem which applies to martingales defined on $[0, L]$, $L \leq \infty$. First recall the fundamental property of conditional expectation that $\|E(f|\mathcal{G})\| \leq E(\|f\|\mathcal{G})$.

Here is a lemma for an optional sampling theorem for the continuous case.

Lemma 31.3.15 *Let $X(t)$ have the property that it is a right continuous nonnegative sub-martingale, $t \geq 0$ such that the filtration $\{\mathcal{F}_t\}$ is normal. Recall this includes $\mathcal{F}_t = \cap_{s>t} \mathcal{F}_s$. Also let τ be a stopping time with values in $[0, T]$. Let $\mathcal{P}_n = \{t_k^n\}_{k=1}^{m_n+1}$ be a sequence of partitions of $[0, T]$ which have the property that*

$$\mathcal{P}_n \subseteq \mathcal{P}_{n+1}, \quad \lim_{n \rightarrow \infty} \|\mathcal{P}_n\| = 0,$$

where $\|\mathcal{P}_n\| \equiv \sup \{|t_k^n - t_{k+1}^n| : k = 1, 2, \dots, m_n\}$. Then let

$$\tau_n(\omega) \equiv \sum_{k=0}^{m_n} t_{k+1}^n \mathcal{X}_{\tau^{-1}((t_k^n, t_{k+1}^n])}(\omega), \quad t_{m_n+1}^n = T$$

It follows that τ_n is a stopping time and also the functions $|X(\tau_n)|$ are uniformly integrable. Furthermore, $|X(\tau)|$ is integrable. Also $X(0), X(\tau_n), X(T)$ is a sub-martingale for the filtration $\mathcal{F}_0, \mathcal{F}_{\tau_n}, \mathcal{F}_T$. If instead $X(t)$ is a martingale having values in a separable Banach space, $X(0), X(\tau_n), X(T)$ is a martingale for the filtration $\mathcal{F}_0, \mathcal{F}_{\tau_n}, \mathcal{F}_T$. In this case, the conclusions about integrability and uniform integrability apply to the sub-martingale $\|X(t)\|$.

Proof: First of all, say $t \in (t_k^n, t_{k+1}^n]$. If $t < t_{k+1}^n$, then

$$[\tau_n \leq t] = [\tau \leq t_k^n] \in \mathcal{F}_{t_k^n} \subseteq \mathcal{F}_t$$

and if $t = t_{k+1}^n$, then

$$[\tau_n \leq t_{k+1}^n] = [\tau \leq t_{k+1}^n] \in \mathcal{F}_{t_{k+1}^n} = \mathcal{F}_t$$

and so τ_n is a stopping time. Thus from Proposition 31.3.11, $X(\tau_n)$ is in $L^1(\Omega)$ the measurability being resolved from Proposition 31.3.8.

Now

$$X(\tau_n) = X\left(\sum_{k=0}^{m_n} t_{k+1}^n \mathcal{X}_{\tau^{-1}((t_k^n, t_{k+1}^n])}(\omega)\right) = \sum_{k=0}^{m_n} X(t_{k+1}^n) \mathcal{X}_{\tau^{-1}((t_k^n, t_{k+1}^n])}(\omega)$$

Now $0, \tau_n, T$ is an increasing list of stopping times. Is it the case that it is a sub-martingale for $\mathcal{F}_0, \mathcal{F}_{\tau_n}, \mathcal{F}_T$?

$$\begin{aligned} E(X(\tau_n) | \mathcal{F}_0) &= \sum_{k=0}^{m_n} E\left(X(t_{k+1}^n) \mathcal{X}_{\tau^{-1}((t_k^n, t_{k+1}^n])} | \mathcal{F}_0\right) \\ &= \sum_{k=0}^{m_n} \mathcal{X}_{\tau^{-1}((t_k^n, t_{k+1}^n])} E(X(t_{k+1}^n) | \mathcal{F}_0) \geq \sum_{k=0}^{m_n} \mathcal{X}_{\tau^{-1}((t_k^n, t_{k+1}^n])} X(0) = X(0) \end{aligned}$$

Now also $X(T) = \sum_{k=0}^{m_n} X(T) \mathcal{X}_{\tau^{-1}((t_k^n, t_{k+1}^n])}(\omega)$ so

$$E(X(T) | \mathcal{F}_{\tau_n}) = \sum_{k=0}^{m_n} E\left(\mathcal{X}_{\tau^{-1}((t_k^n, t_{k+1}^n])} X(T) | \mathcal{F}_{\tau_n}\right). \quad (31.1)$$

What is the value of τ_n on $[\tau \in (t_k^n, t_{k+1}^n]]$? It is t_{k+1}^n , and so from Lemma 31.1.4

$$\begin{aligned} E\left(\mathcal{X}_{\tau^{-1}((t_k^n, t_{k+1}^n])} X(T) | \mathcal{F}_{\tau_n}\right) &= E\left(\mathcal{X}_{\tau^{-1}((t_k^n, t_{k+1}^n])} X(T) | \mathcal{F}_{t_{k+1}^n}\right) \\ &= \mathcal{X}_{\tau^{-1}((t_k^n, t_{k+1}^n])} E(X(T) | \mathcal{F}_{t_{k+1}^n}) \end{aligned}$$

because $\mathcal{X}_{\tau^{-1}((t_k^n, t_{k+1}^n])}$ is $\mathcal{F}_{t_{k+1}^n}$ measurable. Now since X is a sub-martingale,

$$E(X(T) | \mathcal{F}_{t_{k+1}^n}) \geq X(t_{k+1}^n)$$

and so 31.1 implies

$$E(X(T) | \mathcal{F}_{\tau_n}) = \sum_{k=0}^{m_n} E\left(\mathcal{X}_{\tau^{-1}((t_k^n, t_{k+1}^n])} X(T) | \mathcal{F}_{\tau_n}\right)$$

$$\begin{aligned}
&= \sum_{k=0}^{m_n} \mathcal{X}_{\tau^{-1}((t_k^n, t_{k+1}^n])} E(X(T) | \mathcal{F}_{t_{k+1}^n}) \\
&\geq \sum_{k=0}^{m_n} \mathcal{X}_{\tau^{-1}((t_k^n, t_{k+1}^n])} X(t_{k+1}^n) = X(\tau_n)
\end{aligned}$$

Thus this is indeed a sub-martingale. The same argument holds in case X is a martingale. One simply replaces the inequalities with equal signs.

Having shown that $X(0), X(\tau_n), X(T)$ is a sub-martingale,

$$\begin{aligned}
\int_{[X(\tau_n) \geq \lambda]} X(\tau_n) dP &\leq \int_{[X(\tau_n) \geq \lambda]} E(X(T) | \mathcal{F}_{\tau_n}) dP \\
&= \int_{\Omega} E(\mathcal{X}_{[X(\tau_n) \geq \lambda]} X(T) | \mathcal{F}_{\tau_n}) dP \\
&= \int_{[X(\tau_n) \geq \lambda]} X(T) dP
\end{aligned}$$

If the interest were in a martingale where $X(t)$ is in a Banach space, you would simply do all the remaining analysis for the sub-martingale $\|X(t)\|$. Thus, from now on, I will mainly consider a real sub-martingale. From maximal estimates, for example Theorem 29.3.14,

$$P([X(\tau_n) \geq \lambda]) \leq \frac{1}{\lambda} \int_{\Omega} X(T)^+ dP = \frac{1}{\lambda} \int_{\Omega} X(T) dP$$

and now it follows from the above that the random variables $X(\tau_n)$ are equiintegrable. Recall this means that

$$\lim_{\lambda \rightarrow \infty} \sup_n \int_{[X(\tau_n) \geq \lambda]} X(\tau_n) dP = 0$$

Hence they are uniformly integrable and bounded in L^1 .

To verify again that $|X(\tau)|$ is integrable, note that by right continuity, $X(\tau_n) \rightarrow X(\tau)$ pointwise. Apply the Vitali convergence theorem to obtain

$$\int_{\Omega} |X(\tau)| dP = \lim_{n \rightarrow \infty} \int_{\Omega} |X(\tau_n)| dP \leq \int_{\Omega} X(T) dP < \infty. \blacksquare$$

Theorem 31.3.16 *Let $M(t)$ be a right continuous martingale with values in a separable Banach space adapted to a normal filtration. Let σ, τ be two stopping times such that τ is bounded. Then $M(\sigma \wedge \tau) = E(M(\tau) | \mathcal{F}_{\sigma})$. If X is a real sub-martingale, $X(\sigma \wedge \tau) \leq E(X(\tau) | \mathcal{F}_{\sigma})$.*

Proof: Letting $M(t), t \geq 0$ be a martingale with values in a separable Banach space adapted to a filtration \mathcal{F}_t . Let τ_k and σ_k be the discrete stopping times such that τ_k is at least as big as τ but within 2^{-k} of τ discussed earlier. Therefore, from the optional sampling theorem for discrete martingales in Theorem 31.2.1,

$$M(\sigma_n \wedge \tau_n) = E(M(\tau_n) | \mathcal{F}_{\sigma_n})$$

Now let $A \in \mathcal{F}_{\sigma}$. Using Proposition 31.3.8 as needed, $\mathcal{F}_{\sigma} \subseteq \mathcal{F}_{\sigma_n}$ and

$$\int_A M(\sigma_n \wedge \tau_n) dP = \int_A E(M(\tau_n) | \mathcal{F}_{\sigma_n}) dP \stackrel{\mathcal{F}_{\sigma} \subseteq \mathcal{F}_{\sigma_n}}{=} \int_A M(\tau_n) dP$$

Next note that from the right continuity of M , $M(\sigma_n \wedge \tau_n) \rightarrow M(\sigma \wedge \tau)$ and $M(\tau_n) \rightarrow M(\tau)$ and so, by uniform integrability from Proposition 31.3.15, the Vitali convergence theorem applies and we conclude that on passing to a limit,

$$\int_A M(\sigma \wedge \tau) dP = \int_A M(\tau) dP$$

Since $M(\sigma \wedge \tau)$ is $\mathcal{F}_{\sigma \wedge \tau}$ measurable, hence \mathcal{F}_σ measurable, and so, from the definition of conditional expectation, the fact that A is arbitrary implies $M(\sigma \wedge \tau) = E(M(\tau) | \mathcal{F}_\sigma)$.

Now consider the case where X is a sub-martingale. Then by the same observation above about these stopping times and the discrete theory,

$$E(X(\tau_n) | \mathcal{F}_{\sigma_n}) \geq X(\tau_n \wedge \sigma_n)$$

and so, if $A \in \mathcal{F}_\sigma$

$$\int_A X(\tau_n \wedge \sigma_n) dP \leq \int_A E(X(\tau_n) | \mathcal{F}_{\sigma_n}) dP \equiv \int_A X(\tau_n) dP$$

and so, by right continuity and Lemma 31.3.15 and the Vitali convergence theorem, we can pass to a limit and conclude that

$$\int_A X(\tau \wedge \sigma) dP \leq \int_A X(\tau) dP = \int_A E(X(\tau) | \mathcal{F}_\sigma) dP$$

Now $X(\tau \wedge \sigma)$ is $\mathcal{F}_{\tau \wedge \sigma}$ measurable so this function is also \mathcal{F}_σ measurable and so, since the above inequality holds for all $A \in \mathcal{F}_\sigma$, it follows that $X(\tau \wedge \sigma) \leq E(X(\tau) | \mathcal{F}_\sigma)$. ■

Note that a function defined on a countable ordered set such as the integers or equally spaced points is right continuous so the optional sampling theorem for discrete processes is a special case of this one.

31.4 Maximal Inequalities and Stopping Times

As in the case of discrete martingales and sub-martingales, there are maximal inequalities available. Typical ones were presented earlier but here I will use the idea of a stopping time. This gives a typical application of stopping times by making it possible to consider a bounded process and do all the hard work with it and then pass to a limit.

Definition 31.4.1 Let $X(t)$ be a right continuous sub-martingale for $t \in I$ and let $\{\tau_n\}$ be a sequence of stopping times such that $\lim_{n \rightarrow \infty} \tau_n = \infty$. Then X^{τ_n} is called the stopped sub-martingale and it is defined by

$$X^{\tau_n}(t) \equiv X(t \wedge \tau_n).$$

More generally, if τ is a stopping time, the stopped sub-martingale (martingale) is $X^\tau(t) \equiv X(t \wedge \tau)$.

Proposition 31.4.2 The stopped sub-martingale is a sub-martingale.

Proof: By the optional sampling theorem for sub-martingales, Theorem 31.3.16, it follows that for $s < t$,

$$\begin{aligned} E(X^\tau(t) | \mathcal{F}_s) &\equiv E(X(t \wedge \tau) | \mathcal{F}_s) \geq X(t \wedge \tau \wedge s) \\ &= X(\tau \wedge s) \equiv X^\tau(s). \quad \blacksquare \end{aligned}$$

Note that a similar argument would work for martingales.

Theorem 31.4.3 Let $\{X(t)\}$ be a right continuous nonnegative sub-martingale adapted to the normal filtration \mathcal{F}_t for $t \in [0, T]$. Let $p \geq 1$. Define

$$X^*(t) \equiv \sup\{X(s) : 0 < s < t\}, \quad X^*(0) \equiv 0.$$

Then for $\lambda > 0$, if $X(t)^p$ is in $L^1(\Omega)$ for each t ,

$$P([X^*(T) > \lambda]) \leq \frac{1}{\lambda^p} \int \mathcal{X}_{[X^*(T) > \lambda]} X(T)^p dP \quad (31.2)$$

If $X(t)$ is continuous, the above inequality holds without this assumption. In case $p > 1$, and $X(t)$ continuous, then for each $t \leq T$,

$$\left(\int_{\Omega} |X^*(t)|^p dP \right)^{1/p} \leq \frac{p}{p-1} \left(\int_{\Omega} X(T)^p dP \right)^{1/p} \quad (31.3)$$

Proof: The first inequality follows from Theorem 30.5.2. However, it can also be obtained a different way using stopping times. First note that from right continuity, $X^*(t) = \sup\{X(d) : d \in D\}$ where D is a dense countable set in $(0, t)$. Therefore, $X^*(t)$ is always measurable.

First I will assume $X(t)$ is a bounded sub-martingale. These certainly exist. Just take a bounded stopping time τ and consider X^τ .

Define the stopping time

$$\tau \equiv \inf\{t > 0 : X(t) > \lambda\} \wedge T.$$

(The infimum over an empty set will equal ∞ .) This is a stopping time by 31.3.9 because it is just a continuous function of the first hitting time of an open set. Also from the definition of X^* in which the supremum is taken over an open interval,

$$[\tau < t] = [X^*(t) > \lambda]$$

Note this also shows $X^*(t)$ is \mathcal{F}_t measurable. Then it follows that $X^p(t)$ is also a sub-martingale since r^p is increasing and convex. By the optional sampling theorem,

$$X(0)^p, X(\tau)^p, X(T)^p$$

is a sub-martingale. Recall $X(\sigma \wedge \tau) \leq E(X(\tau) | \mathcal{F}_\sigma)$ when τ is bounded. I need to verify that

$$E(X(T)^p | \mathcal{F}_\tau) \geq X(\tau)^p, E(X(\tau)^p | \mathcal{F}_0) \geq X(0)^p.$$

But from the optional sampling theorem Theorem 31.3.16

$$\begin{aligned} E(X(T)^p | \mathcal{F}_\tau) &\geq X(T \wedge \tau)^p = X(\tau)^p \\ E(X(\tau)^p | \mathcal{F}_0) &\geq X(\tau \wedge 0)^p = X(0)^p \end{aligned}$$

Also $[\tau < T] \in \mathcal{F}_\tau$. Recall that A is \mathcal{F}_τ measurable means $A \cap [\tau \leq t] \in \mathcal{F}_t$. Since τ is a stopping time, $[\tau \leq T] \cap [\tau \leq t] = [\tau \leq t] \in \mathcal{F}_t$ and so $[\tau \leq T] \in \mathcal{F}_\tau$.

$$\int_{[\tau < T]} X(\tau)^p dP \leq \int_{[\tau < T]} E(X(T)^p | \mathcal{F}_\tau) dP = \int_{[\tau < T]} X(T)^p dP$$

By right continuity, on $[\tau < T]$, $X(\tau) \geq \lambda$. Therefore,

$$\begin{aligned} \lambda^p P([X^*(T) > \lambda]) &= \lambda^p P([\tau < T]) \\ &\leq \int_{[\tau < T]} X(\tau)^p dP = \int_{[\tau < T]} E(X(T)^p | \mathcal{F}_\tau) dP = \int_{[X^*(T) > \lambda]} X(T)^p dP \end{aligned}$$

This proves 31.2 in case X is bounded. In general case, suppose X is not just right continuous but also continuous.

Next let $\{\tau_n\}$ be a “localizing sequence” given by

$$\tau_n \equiv \inf \{t : X(t) > n\}.$$

If $t < \tau_n$, then $X(t) \leq n$ by definition of τ_n . Could $X(\tau_n) > n$? If so, then by continuity, $X(t) > n$ for some $t < \tau_n$ so τ_n was not chosen correctly. Thus X^{τ_n} is bounded because $X(\tau_n \wedge t) \leq n$, and so from what was just shown,

$$\lambda^p P([(X^{\tau_n})^*(T) > \lambda]) \leq \int_{[(X^{\tau_n})^*(T) > \lambda]} (X^{\tau_n}(T))^p dP$$

Then $(X^{\tau_n})(T)$ is increasing to $X(T)$ and $(X^{\tau_n})^*(T)$ increases to $X^*(T)$ as $n \rightarrow \infty$ so 31.2 follows from the monotone convergence theorem. This proves 31.2.

Let X^{τ_n} be as just defined. Thus it is a bounded sub-martingale. To save on notation, the X in the following argument is really X^{τ_n} . This is done so that all the integrals are finite. If $p > 1$, then from the first part,

$$\begin{aligned} \int_{\Omega} |X^*(t)|^p dP &\leq \int_{\Omega} |X^*(T)|^p dP = \int_0^\infty p \lambda^{p-1} \overbrace{P([X^*(T) > \lambda])}^{\leq \frac{1}{\lambda} \int \mathcal{X}_{[X^*(T) > \lambda]} X(T) dP} d\lambda \\ &\leq p \int_0^\infty \lambda^{p-1} \frac{1}{\lambda} \int_{\mathcal{X}_{[X^*(T) > \lambda]} X(T) dP d\lambda \end{aligned}$$

By Lemma 29.3.13, applied to the second half of the above and using Holder’s inequality,

$$\begin{aligned} \int_{\Omega} |X^*(T)|^p dP &\leq p \int_{\Omega} X(T) \int_0^{X^*(T)} \lambda^{p-2} d\lambda dP = p \int_{\Omega} X(T) \frac{X^*(T)^{p-1}}{p-1} dP \\ &\leq \frac{p}{p-1} \left(\int_{\Omega} X^*(T)^p dP \right)^{1/p'} \left(\int_{\Omega} X(T)^p dP \right)^{1/p} \end{aligned}$$

Now divide both sides by $(\int_{\Omega} X^*(T)^p dP)^{1/p'}$ and restore X^{τ_n} for X .

$$\left(\int_{\Omega} X^{\tau_n*}(T)^p dP \right)^{1/p} \leq \frac{p}{p-1} \left(\int_{\Omega} X^{\tau_n}(T)^p dP \right)^{1/p}$$

Now let $n \rightarrow \infty$ and use the monotone convergence theorem to obtain the inequality of the theorem 31.3. ■

Here is another sort of maximal inequality in which $X(t)$ is not assumed nonnegative. A version of this was also presented earlier.

Theorem 31.4.4 Let $\{X(t)\}$ be a right continuous sub-martingale adapted to the normal filtration \mathcal{F}_t for $t \in [0, T]$ and $X^*(t)$ defined as in Theorem 31.4.3

$$X^*(t) \equiv \sup \{X(s) : 0 < s < t\}, \quad X^*(0) \equiv 0,$$

$$P([X^*(T) > \lambda]) \leq \frac{1}{\lambda} E(|X(T)|) \quad (31.4)$$

For $t > 0$, let

$$X_*(t) = \inf \{X(s) : s < t\}.$$

Then

$$P([X_*(T) < -\lambda]) \leq \frac{1}{\lambda} E(|X(T)| + |X(0)|) \quad (31.5)$$

Also

$$P([\sup \{|X(s)| : s < T\} > \lambda])$$

$$\leq \frac{2}{\lambda} E(|X(T)| + |X(0)|) \quad (31.6)$$

Proof: The function $f(r) = r^+ \equiv \frac{1}{2}(|r| + r)$ is convex and increasing. Therefore, $X^+(t)$ is also a sub-martingale but this one is nonnegative. Also

$$[X^*(T) > \lambda] = [(X^+)^*(T) > \lambda]$$

and so from Theorem 31.4.3,

$$P([X^*(T) > \lambda]) = P([(X^+)^*(T) > \lambda]) \leq \frac{1}{\lambda} E(X^+(T)) \leq \frac{1}{\lambda} E(|X(T)|).$$

Next let

$$\tau = \min(\inf \{t : X(t) < -\lambda\}, T)$$

then as before, $X(0), X(\tau), X(T)$ is a sub-martingale and so

$$\int_{[\tau < T]} X(\tau) dP + \int_{[\tau = T]} X(\tau) dP = \int_{\Omega} X(\tau) dP \geq \int_{\Omega} X(0) dP$$

Now for $\omega \in [\tau < T]$, $X(t)(\omega) < -\lambda$ for some $t < T$ and so it follows that by right continuity, $X(\tau)(\omega) \leq -\lambda$. therefore,

$$-\lambda \int_{[\tau < T]} dP \geq - \int_{[\tau = T]} X(T) dP + \int_{\Omega} X(0) dP$$

If $X_*(T) < -\lambda$, then from the definition given above, there exists $t < T$ such that $X(t) < -\lambda$ and so $\tau < T$. If $\tau < T$, then by definition, there exists $t < T$ such that $X(t) < -\lambda$ and so $X_*(T) < -\lambda$. Hence $[\tau < T] = [X_*(T) < -\lambda]$. It follows that

$$P([X_*(T) < -\lambda]) = P([\tau < T])$$

$$\leq \frac{1}{\lambda} \int_{[\tau = T]} X(T) dP - \frac{1}{\lambda} \int_{\Omega} X(0) dP \leq \frac{1}{\lambda} E(|X(T)| + |X(0)|)$$

and this proves 31.5.

Finally, combining the above two inequalities,

$$P([\sup \{|X(s)| : s < T\} > \lambda]) = P([X_*(T) < -\lambda]) + P([X^*(T) > \lambda])$$

$$\leq \frac{2}{\lambda} E(|X(T)| + |X(0)|). \blacksquare$$

31.5 Continuous Sub-martingale Convergence

Here, $\{Y(t)\}$ will be a continuous sub-martingale and $a < b$. Let $X(t) \equiv (Y(t) - a)_+ + a$ so $X(0) \geq a$. Then X is also a sub-martingale. It is an increasing convex function of one. If $Y(t)$ has an upcrossing of $[a, b]$, then $X(t)$ starts off at a and ends up at least as large as b . If $X(t)$ has an upcrossing of $[a, b]$ then it must start off at a since it cannot be smaller and it ends up at least as large as b . Thus we can count the upcrossings of $Y(t)$ by considering the upcrossings of $X(t)$ and $X(t)$ is always at least as large as a .

The next task is to consider an upcrossing estimate as was done before for discrete sub-martingales.

$$\begin{aligned}\tau_0 &\equiv \min(\inf\{t > 0 : X(t) = a\}, M), \\ \tau_1 &\equiv \min(\inf\{t > 0 : (X(t \vee \tau_0) - X(\tau_0))_+ = b - a\}, M), \\ \tau_2 &\equiv \min(\inf\{t > 0 : (X(\tau_1) - X(t \vee \tau_1))_+ = b - a\}, M), \\ \tau_3 &\equiv \min(\inf\{t > 0 : (X(t \vee \tau_2) - X(\tau_2))_+ = b - a\}, M), \\ \tau_4 &\equiv \min(\inf\{t > 0 : (X(\tau_3) - X(t \vee \tau_3))_+ = b - a\}, M), \\ &\vdots\end{aligned}$$

If $X(t)$ is never a , then $\tau_0 \equiv M$ where we assume $t \in [0, M]$ and there are no upcrossings. It is obvious $\tau_1 \geq \tau_0$ since otherwise, the inequality could not hold. Thus the evens have $X(\tau_{2k}) = a$ and $X(\tau_{2k+1}) = b$. The following lemma follows from Corollary 31.3.14.

Lemma 31.5.1 *The above τ_i are stopping times for $t \in [0, M]$.*

Note that in the above, if $\eta = M$, then $\sigma = M$ also. Thus in the definition of the τ_i , if any $\tau_i = M$, it follows that also $\tau_{i+1} = M$ and so there is no change in the stopping times. Also note that these stopping times τ_i are increasing as i increases.

Let

$$U_{[a,b]}^{nM} \equiv \lim_{\varepsilon \rightarrow 0} \sum_{k=0}^n \frac{X(\tau_{2k+1}) - X(\tau_{2k})}{\varepsilon + X(\tau_{2k+1}) - X(\tau_{2k})}$$

Note that if an upcrossing occurs after τ_{2k} on $[0, M]$, then $\tau_{2k+1} > \tau_{2k}$ because there exists t such that

$$(X(t \vee \tau_{2k}) - X(\tau_{2k}))_+ = b - a$$

However, you could have $\tau_{2k+1} > \tau_{2k}$ without an upcrossing occurring. This happens when $\tau_{2k} < M$ and $\tau_{2k+1} = M$ which may mean that $X(t)$ never again climbs to b . You break the sum into those terms where $X(\tau_{2k+1}) - X(\tau_{2k}) = b - a$ and those where this is less than $b - a$. Suppose for a fixed ω , the terms where the difference is $b - a$ are for $k \leq m$. Then there might be a last term for which $X(\tau_{2k+1}) - X(\tau_{2k}) < b - a$ because it fails to complete the up crossing. There is only one of these at $k = m + 1$. Then the above sum is

$$\begin{aligned}&\leq \frac{1}{b-a} \sum_{k=0}^m X(\tau_{2k+1}) - X(\tau_{2k}) + \frac{X(M) - a}{\varepsilon + X(M) - a} \\&\leq \frac{1}{b-a} \sum_{k=0}^n X(\tau_{2k+1}) - X(\tau_{2k}) + \frac{X(M) - a}{\varepsilon + X(M) - a} \\&\leq \frac{1}{b-a} \sum_{k=0}^n X(\tau_{2k+1}) - X(\tau_{2k}) + 1\end{aligned}$$

Then $U_{[a,b]}^{nM}$ is clearly a random variable which is at least as large as the number of upcrossings occurring for $t \leq M$ using only $2n+1$ of the stopping times. From the optional sampling theorem,

$$\begin{aligned} E(X(\tau_{2k}) - X(\tau_{2k-1})) &= \int_{\Omega} X(\tau_{2k}) - X(\tau_{2k-1}) dP \\ &= \int_{\Omega} E(X(\tau_{2k}) | \mathcal{F}_{\tau_{2k-1}}) - X(\tau_{2k-1}) dP \\ &\geq \int_{\Omega} X(\tau_{2k-1}) - X(\tau_{2k-1}) dP = 0 \end{aligned}$$

Note that $X(\tau_{2k}) = a$ while $X(\tau_{2k-1}) = b$ so the above may seem surprising. However, the two stopping times can both equal M so this is actually possible. For example, it could happen that $X(t) = a$ for all $t \in [0, M]$.

Next, take the expectation of both sides,

$$\begin{aligned} E(U_{[a,b]}^{nM}) &\leq \frac{1}{b-a} \sum_{k=0}^n E(X(\tau_{2k+1})) - E(X(\tau_{2k})) + 1 \\ &\leq \frac{1}{b-a} \sum_{k=0}^n E(X(\tau_{2k+1})) - E(X(\tau_{2k})) + \frac{1}{b-a} \sum_{k=1}^n E(X(\tau_{2k})) - E(X(\tau_{2k-1})) + 1 \\ &= \frac{1}{b-a} (E(X(\tau_1)) - E(X(\tau_0))) + \frac{1}{b-a} \sum_{k=1}^n E(X(\tau_{2k+1})) - E(X(\tau_{2k-1})) + 1 \\ &\leq \frac{1}{b-a} (E(X(\tau_{2n+1})) - E(X(\tau_0))) + 1 \\ &\leq \frac{1}{b-a} (E(X(M)) - a) + 1 \end{aligned}$$

which does not depend on n . The last inequality follows because $0 \leq \tau_{2n+1} \leq M$ and $X(t)$ is a sub-martingale. Let $n \rightarrow \infty$ to obtain

$$E(U_{[a,b]}^M) \leq \frac{1}{b-a} (E(X(M)) - a) + 1$$

where $U_{[a,b]}^M$ is an upper bound to the number of upcrossings of $\{X(t)\}$ on $[0, M]$. This proves the following interesting upcrossing estimate.

Lemma 31.5.2 *Let $\{Y(t)\}$ be a continuous sub-martingale adapted to a normal filtration \mathcal{F}_t for $t \in [0, M]$. Then if $U_{[a,b]}^M$ is defined as the above upper bound to the number of upcrossings of $\{Y(t)\}$ for $t \in [0, M]$, then this is a random variable and*

$$\begin{aligned} E(U_{[a,b]}^M) &\leq \frac{1}{b-a} (E(Y(M) - a)_+ + a - a) + 1 \\ &= \frac{1}{b-a} E(|Y(M)|) + \frac{1}{b-a} |a| + 1 \end{aligned}$$

With this it is easy to prove a continuous sub-martingale convergence theorem.

Theorem 31.5.3 *Let $\{X(t)\}$ be a continuous sub-martingale adapted to a normal filtration such that*

$$\sup_t \{E(|X(t)|)\} = C < \infty.$$

Then there exists $X_\infty \in L^1(\Omega)$ such that

$$\lim_{t \rightarrow \infty} X(t)(\omega) = X_\infty(\omega) \text{ a.e. } \omega.$$

Proof: Let $U_{[a,b]}$ be defined by

$$U_{[a,b]} = \lim_{M \rightarrow \infty} U_{[a,b]}^M.$$

Thus the random variable $U_{[a,b]}$ is an upper bound for the number of upcrossings. From Lemma 31.5.2 and the assumption of this theorem, there exists a constant C independent of M such that

$$E(U_{[a,b]}^M) \leq \frac{C}{b-a} + 1.$$

Letting $M \rightarrow \infty$, it follows from monotone convergence theorem that

$$E(U_{[a,b]}) \leq \frac{C}{b-a} + 1$$

also. Therefore, there exists a set of measure 0 N_{ab} such that if $\omega \notin N_{ab}$, then $U_{[a,b]}(\omega) < \infty$. That is, there are only finitely many upcrossings. Now let

$$N = \cup \{N_{ab} : a, b \in \mathbb{Q}\}.$$

It follows that for $\omega \notin N$, it cannot happen that

$$\limsup_{t \rightarrow \infty} X(t)(\omega) - \liminf_{t \rightarrow \infty} X(t)(\omega) > 0$$

because if this expression is positive, there would be arbitrarily large values of t where $X(t)(\omega) > b$ and arbitrarily large values of t where $X(t)(\omega) < a$ where a, b are rational numbers chosen such that

$$\limsup_{t \rightarrow \infty} X(t)(\omega) > b > a > \liminf_{t \rightarrow \infty} X(t)(\omega)$$

Thus there would be infinitely many upcrossings which is not allowed for $\omega \notin N$. Therefore, the limit $\lim_{t \rightarrow \infty} X(t)(\omega)$ exists for a.e. ω . Let $X_\infty(\omega)$ equal this limit for $\omega \notin N$ and let $X_\infty(\omega) = 0$ for $\omega \in N$. Then X_∞ is measurable and by Fatou's lemma,

$$\int_{\Omega} |X_\infty(\omega)| dP \leq \liminf_{n \rightarrow \infty} \int_{\Omega} |X(n)(\omega)| dP < C. \blacksquare$$

Now here is an interesting result of Doob.

Theorem 31.5.4 *Let $\{M(t)\}$ be a continuous real martingale adapted to the normal filtration \mathcal{F}_t . Then the following are equivalent.*

1. *The random variables $M(t)$ are equi-integrable.*

2. There exists $M(\infty) \in L^1(\Omega)$ such that $\lim_{t \rightarrow \infty} \|M(\infty) - M(t)\|_{L^1(\Omega)} = 0$.

In this case, $M(t) = E(M(\infty) | \mathcal{F}_t)$ and convergence also takes place pointwise.

Proof: Suppose the equi-integrable condition. Then there exists λ large enough that for all t ,

$$\int_{[|M(t)| \geq \lambda]} |M(t)| dt < 1.$$

It follows that for all t ,

$$\begin{aligned} \int_{\Omega} |M(t)| dP &= \int_{[|M(t)| \geq \lambda]} |M(t)| dP + \int_{[|M(t)| < \lambda]} |M(t)| dP \\ &\leq 1 + \lambda. \end{aligned}$$

Since the martingale is bounded in L^1 , by Theorem 31.5.3 there exists $M(\infty) \in L^1(\Omega)$ such that $\lim_{t \rightarrow \infty} M(t)(\omega) = M(\infty)(\omega)$ pointwise a.e. By the assumption $\{M(t)\}$ are equi-integrable, it follows from Proposition 10.9.6 these functions are uniformly integrable. Then by the Vitali convergence theorem, Theorem 10.9.7, if $t_n \rightarrow \infty$, then

$$\|M(t_n) - M(\infty)\|_{L^1(\Omega)} \rightarrow 0$$

Next suppose there is a function $M(\infty)$ to which $M(t)$ converges in $L^1(\Omega)$. Then for t fixed and $A \in \mathcal{F}_t$, then as $s \rightarrow \infty, s > t$

$$\begin{aligned} \int_A M(t) dP &= \int_A E(M(s) | \mathcal{F}_t) dP \equiv \int_A M(s) dP \\ &\rightarrow \int_A M(\infty) dP = \int_A E(M(\infty) | \mathcal{F}_t) dP \end{aligned}$$

which shows $E(M(\infty) | \mathcal{F}_t) = M(t)$ a.e. since $A \in \mathcal{F}_t$ is arbitrary. By Theorem 24.12.1,

$$\begin{aligned} \int_{[|M(t)| \geq \lambda]} |M(t)| dP &= \int_{[|M(t)| \geq \lambda]} |E(M(\infty) | \mathcal{F}_t)| dP \\ &\leq \int_{[|M(t)| \geq \lambda]} E(|M(\infty)| | \mathcal{F}_t) dP \\ &= \int_{[|M(t)| \geq \lambda]} |M(\infty)| dP \end{aligned} \tag{31.7}$$

Now from this,

$$\begin{aligned} \lambda P([|M(t)| \geq \lambda]) &\leq \int_{[|M(t)| \geq \lambda]} |M(t)| dP \leq \int_{\Omega} |E(M(\infty) | \mathcal{F}_t)| dP \\ &\leq \int_{\Omega} E(|M(\infty)| | \mathcal{F}_t) dP = \int_{\Omega} |M(\infty)| dP \end{aligned}$$

and so

$$P([|M(t)| \geq \lambda]) \leq \frac{C}{\lambda}$$

From 31.7, this shows $\{M(t)\}$ is equi-integrable hence uniformly integrable because this is true of the single function $|M(\infty)|$. ■

31.6 Hitting This Before That

Let $\{M(t)\}$ be a real valued continuous martingale for $t \in [0, T]$ where $T \leq \infty$ and $M(0) = 0$. In case $T = \infty$, assume the conditions of Theorem 31.5.4 are satisfied. Thus, according to these conditions, there exists $M(\infty)$ and the $M(t)$ are equi-integrable. With the Doob optional sampling theorem it is possible to estimate the probability that $M(t)$ hits a before it hits b where $a < 0 < b$. There is no loss of generality in assuming $T = \infty$ since if it is less than ∞ , you could just let $M(t) \equiv M(T)$ for all $t > T$. In this case, the equiintegrability of the $M(t)$ follows because for $t < T$,

$$\begin{aligned} \int_{[|M(t)| > \lambda]} |M(t)| dP &= \int_{[|M(t)| > \lambda]} |E(M(T) | \mathcal{F}_t)| dP \\ &\leq \int_{[|M(t)| > \lambda]} E(|M(T)| | \mathcal{F}_t) dP = \int_{[|M(t)| > \lambda]} |M(T)| dP \end{aligned}$$

and from Theorem 31.4.4,

$$P(|M(t)| > \lambda) \leq P([M^*(t) > \lambda]) \leq \frac{1}{\lambda} \int_{\Omega} |M(T)| dP.$$

Definition 31.6.1 Let M be a process adapted to the filtration \mathcal{F}_t and let τ be a stopping time. Then M^τ , called the stopped process is defined by

$$M^\tau(t) \equiv M(\tau \wedge t).$$

With this definition, here is a simple lemma. I will use this lemma whenever convenient without comment.

Lemma 31.6.2 Let M be a right continuous martingale adapted to the normal filtration \mathcal{F}_t and let τ be a stopping time. Then M^τ is also a martingale adapted to the filtration \mathcal{F}_t . The same is true for a sub-martingale.

Proof: Let $s < t$. By the Doob optional sampling theorem,

$$E(M^\tau(t) | \mathcal{F}_s) \equiv E(M(\tau \wedge t) | \mathcal{F}_s) = M(\tau \wedge t \wedge s) = M^\tau(s).$$

As for a sub-martingale $X(t)$, for $s < t$

$$E(X^\tau(t) | \mathcal{F}_s) \equiv E(X(\tau \wedge t) | \mathcal{F}_s) \geq X(\tau \wedge t \wedge s) \equiv X^\tau(s). \blacksquare$$

Theorem 31.6.3 Let $\{M(t)\}$ be a continuous real valued martingale adapted to the normal filtration \mathcal{F}_t and let

$$M^* \equiv \sup\{|M(t)| : t \geq 0\}$$

and $M(0) = 0$. Letting

$$\tau_x \equiv \inf\{t > 0 : M(t) = x\}$$

Then if $a < 0 < b$ the following inequalities hold.

$$(b-a)P([\tau_b \leq \tau_a]) \geq -aP([M^* > 0]) \geq (b-a)P([\tau_b < \tau_a])$$

and

$$(b-a)P([\tau_a < \tau_b]) \leq bP([M^* > 0]) \leq (b-a)P([\tau_a \leq \tau_b]).$$

In words, $P([\tau_b \leq \tau_a])$ is the probability that $M(t)$ hits b no later than when it hits a . (Note that if $\tau_a = \infty = \tau_b$ then you would have $[\tau_a = \tau_b]$.)

Proof: For $x \in \mathbb{R}$, define

$$\tau_x \equiv \inf\{t \in \mathbb{R} \text{ such that } M(t) = x\}$$

with the usual convention that $\inf(\emptyset) = \infty$. Let $a < 0 < b$ and let

$$\tau = \tau_a \wedge \tau_b$$

Then the following claim will be important.

Claim: $E(M(\tau)) = 0$.

Proof of the claim: Let $t > 0$. Then by the Doob optional sampling theorem,

$$E(M(\tau \wedge t)) = E(E(M(t) | \mathcal{F}_\tau)) = E(M(t)) \quad (31.8)$$

$$= E(E(M(t) | \mathcal{F}_0)) = E(M(0)) = 0. \quad (31.9)$$

Observe the martingale M^τ must be bounded because it is stopped when $M(t)$ equals either a or b . There are two cases according to whether $\tau = \infty$. If $\tau = \infty$, then $M(t)$ never hits a or b so $M(t)$ has values between a and b . In this case $M^\tau(t) = M(t) \in [a, b]$. On the other hand, you could have $\tau < \infty$. Then in this case $M^\tau(t)$ is eventually equal to either a or b depending on which it hits first. In either case, the martingale M^τ is bounded and by the martingale convergence theorem, Theorem 31.5.3, there exists $M^\tau(\infty)$ such that

$$\lim_{t \rightarrow \infty} M^\tau(t)(\omega) = M^\tau(\infty)(\omega) = M(\tau)(\omega)$$

and since the $M^\tau(t)$ are bounded, the dominated convergence theorem implies

$$E(M(\tau)) = \lim_{t \rightarrow \infty} E(M(\tau \wedge t)) = 0.$$

This proves the claim.

Let

$$M^*(\omega) \equiv \sup\{|M(t)(\omega)| : t \in [0, \infty]\}.$$

Also note that $[\tau_a = \tau_b] = [\tau = \infty]$. This is because $a \neq b$. If $M(t) = a$, then $M(t) \neq b$ so it cannot happen that these are equal at any finite time. But if $\tau = \infty$, then both $\tau_a, \tau_b = \infty$. Now from the claim,

$$\begin{aligned} 0 &= E(M(\tau)) = \int_{[\tau_a < \tau_b]} M(\tau) dP \\ &\quad + \int_{[\tau_b < \tau_a]} M(\tau) dP + \int_{[\tau_a = \tau_b] \cap [M^* > 0]} M(\infty) dP \\ &\quad + \int_{[\tau_a = \tau_b] \cap [M^* = 0]} M(\infty) dP \end{aligned} \quad (31.10)$$

The last term equals 0. By continuity, $M(\tau)$ is either equal to a or b depending on whether $\tau_a < \tau_b$ or $\tau_b < \tau_a$. Thus

$$\begin{aligned} 0 &= E(M(\tau)) = aP([\tau_a < \tau_b]) \\ &\quad + bP([\tau_b < \tau_a]) + \int_{[\tau_a = \tau_b] \cap [M^* > 0]} M(\infty) dP \end{aligned} \quad (31.11)$$

Consider this last term. By the definition, $[\tau_a = \tau_b]$ corresponds to $M(t)$ never hitting either a or b . Since $M(0) = 0$, this can only happen if $M(t)$ has values in $[a, b]$. Therefore, this last term satisfies

$$\begin{aligned} aP([\tau_a = \tau_b] \cap [M^* > 0]) &\leq \int_{[\tau_a = \tau_b] \cap [M^* > 0]} M(\infty) dP \\ &\leq bP([\tau_a = \tau_b] \cap [M^* > 0]) \end{aligned} \quad (31.12)$$

Obviously the following inequality holds because on the left you have

$$aP([\tau_a = \tau_b] \cap [M^* > 0])$$

and on the right you have the larger $bP([\tau_a = \tau_b] \cap [M^* > 0])$. That 0 is in the middle follows from 31.11.

$$\begin{aligned} aP([\tau_a = \tau_b] \cap [M^* > 0]) + aP([\tau_a < \tau_b]) + bP([\tau_b < \tau_a]) &\leq \\ 0 \leq bP([\tau_a = \tau_b] \cap [M^* > 0]) + aP([\tau_a < \tau_b]) + bP([\tau_b < \tau_a]) \end{aligned} \quad (31.13)$$

Note that $[\tau_b < \tau_a], [\tau_a < \tau_b] \subseteq [M^* > 0]$ and so

$$[\tau_b < \tau_a] \cup [\tau_a < \tau_b] \cup ([\tau_a = \tau_b] \cap [M^* > 0]) = [M^* > 0] \quad (31.14)$$

The following diagram may help in keeping track of the various substitutions.

$[\tau_a < \tau_b]$	$[\tau_b < \tau_a]$	$[\tau_b = \tau_a] \cap [M^* > 0]$
---------------------	---------------------	------------------------------------

Left side of 31.13

From 31.14, this yields on substituting for $P([\tau_a < \tau_b])$

$$\begin{aligned} 0 &\geq aP([\tau_a = \tau_b] \cap [M^* > 0]) + a[P([M^* > 0]) - P([\tau_a \geq \tau_b] \cap [M^* > 0])] \\ &\quad + bP([\tau_b < \tau_a]) \end{aligned}$$

and so since $[\tau_a \neq \tau_b] \subseteq [M^* > 0]$,

$$0 \geq a[P([M^* > 0]) - P([\tau_a > \tau_b])] + bP([\tau_b < \tau_a])$$

$$\boxed{-aP([M^* > 0]) \geq (b-a)P([\tau_b < \tau_a])} \quad (31.15)$$

Next use 31.14 to substitute for $P([\tau_b < \tau_a])$

$$\begin{aligned} 0 &\geq aP([\tau_a = \tau_b] \cap [M^* > 0]) + aP([\tau_a < \tau_b]) + bP([\tau_b < \tau_a]) \\ &= aP([\tau_a = \tau_b] \cap [M^* > 0]) + aP([\tau_a < \tau_b]) \\ &\quad + b[P([M^* > 0]) - P([\tau_a \leq \tau_b] \cap [M^* > 0])] \\ &= aP([\tau_a \leq \tau_b] \cap [M^* > 0]) + b[P([M^* > 0]) - P([\tau_a \leq \tau_b] \cap [M^* > 0])] \end{aligned}$$

and so

$$\boxed{(b-a)P([\tau_a \leq \tau_b]) \geq bP([M^* > 0])} \quad (31.16)$$

Right side of 31.13

From 31.14, used to substitute for $P([\tau_a < \tau_b])$ this yields

$$\begin{aligned}
 0 &\leq bP([\tau_a = \tau_b] \cap [M^* > 0]) + aP([\tau_a < \tau_b]) + bP([\tau_b < \tau_a]) \\
 &= bP([\tau_a = \tau_b] \cap [M^* > 0]) + a[P([M^* > 0]) - P([\tau_a \geq \tau_b] \cap [M^* > 0])] \\
 &\quad + bP([\tau_b < \tau_a]) \\
 &= bP([\tau_a \geq \tau_b] \cap [M^* > 0]) + a[P([M^* > 0]) - P([\tau_a \geq \tau_b] \cap [M^* > 0])]
 \end{aligned}$$

and so

$$(b-a)P([\tau_a \geq \tau_b]) \geq -aP([M^* > 0]) \quad (31.17)$$

Next use 31.14 to substitute for the term $P([\tau_b < \tau_a])$ and write

$$\begin{aligned}
 0 &\leq bP([\tau_a = \tau_b] \cap [M^* > 0]) + aP([\tau_a < \tau_b]) + bP([\tau_b < \tau_a]) \\
 &= bP([\tau_a = \tau_b] \cap [M^* > 0]) + aP([\tau_a < \tau_b]) \\
 &\quad + b[P([M^* > 0]) - P([\tau_a \leq \tau_b] \cap [M^* > 0])] \\
 &= aP([\tau_a < \tau_b]) + bP([M^* > 0]) - bP([\tau_a < \tau_b] \cap [M^* > 0]) \\
 &= aP([\tau_a < \tau_b]) + bP([M^* > 0]) - bP([\tau_a < \tau_b])
 \end{aligned}$$

and so

$$(b-a)P([\tau_a < \tau_b]) \leq bP([M^* > 0]) \quad (31.18)$$

Now the boxed in formulas in 31.15 - 31.18 yield the conclusion of the theorem. ■

Note $P([\tau_a < \tau_b])$ means $M(t)$ hits a before it hits b with other occurrences of similar expressions being defined similarly.

31.7 The Space $\mathcal{M}_T^p(E)$

Here $p \geq 1$. Also, we assume the filtration is a normal filtration.

Definition 31.7.1 Then $M \in \mathcal{M}_T^p(E)$ if $t \rightarrow M(t)(\omega)$ is continuous for a.e. ω and $M(t)$ is adapted, and

$$E \left(\sup_{t \in [0, T]} \|M(t)\|^p \right) < \infty$$

Here E is a separable Banach space.

Proposition 31.7.2 Define a norm on $\mathcal{M}_T^p(E)$ by

$$\|M\|_{\mathcal{M}_T^p(E)} \equiv E \left(\sup_{t \in [0, T]} \|M(t)\|^p \right)^{1/p}.$$

Then with this norm, $\mathcal{M}_T^p(E)$ is a Banach space. Also, a Cauchy sequence in this space has a subsequence which converges uniformly for all ω off a set of measure zero. Those M in $\mathcal{M}_T^p(E)$ which are martingales constitute a closed subspace of $\mathcal{M}_T^p(E)$. If σ is a stopping time, then if $M \in \mathcal{M}_T^p(E)$, so is M^σ and $\|M\|_{\mathcal{M}_T^p(E)} \geq \|M^\sigma\|_{\mathcal{M}_T^p(E)}$. Thus if $M_n \rightarrow M$ in $\mathcal{M}_T^p(E)$, then $M_n^\sigma \rightarrow M^\sigma$ in $\mathcal{M}_T^p(E)$.

Proof: First it is good to observe that $\sup_{t \in [0, T]} \|M(t)\|^p$ is measurable. This follows because of the continuity of $t \rightarrow M(t)$. Let D be a dense countable set in $[0, T]$. Then by continuity,

$$\sup_{t \in [0, T]} \|M(t)\|^p = \sup_{t \in D} \|M(t)\|^p$$

and the expression on the right is measurable because D is countable.

Next it is necessary to show this is a norm. It is clear that $\|M\|_{\mathcal{M}_T^p(E)} \geq 0$ and equals 0 only if $0 = E \left(\sup_{t \in [0, T]} \|M(t)\|^p \right)$ which requires $M(t) = 0$ for all t for ω off a set of measure zero so that $M = 0$. It is also clear that $\|\alpha M\|_{\mathcal{M}_T^p(E)} = |\alpha| \|M\|_{\mathcal{M}_T^p(E)}$. It remains to check the triangle inequality. Let $M, N \in \mathcal{M}_T^p(E)$.

$$\begin{aligned} \|M+N\|_{\mathcal{M}_T^p(E)} &\equiv E \left(\sup_{t \in [0, T]} \|M(t) + N(t)\|^p \right)^{1/p} \\ &\leq E \left(\sup_{t \in [0, T]} (\|M(t)\| + \|N(t)\|)^p \right)^{1/p} \\ &\leq E \left(\left(\sup_{t \in [0, T]} \|M(t)\| + \sup_{t \in [0, T]} \|N(t)\| \right)^p \right)^{1/p} \\ &\equiv \left(\int_{\Omega} \left(\sup_{t \in [0, T]} \|M(t)\| + \sup_{t \in [0, T]} \|N(t)\| \right)^p dP \right)^{1/p} \\ &\leq \left(\int_{\Omega} \left(\sup_{t \in [0, T]} \|M(t)\| \right)^p dP \right)^{1/p} + \left(\int_{\Omega} \left(\sup_{t \in [0, T]} \|N(t)\| \right)^p dP \right)^{1/p} \\ &\equiv \|M\|_{\mathcal{M}_T^p(E)} + \|N\|_{\mathcal{M}_T^p(E)} \end{aligned}$$

Next consider the claim that $\mathcal{M}_T^p(E)$ is a Banach space. Let $\{M_n\}$ be a Cauchy sequence. Then

$$E \left(\sup_{t \in [0, T]} \|M_n(t) - M_m(t)\|^p \right) \rightarrow 0 \quad (31.19)$$

as $m, n \rightarrow \infty$. Now

$$P \left(\sup_{t \in [0, T]} \|M_n(t) - M_m(t)\| > \lambda \right) \leq \frac{1}{\lambda^p} E \left(\sup_{t \in [0, T]} \|M_n(t) - M_m(t)\|^p \right)$$

Therefore, one can extract a subsequence $\{M_{n_k}\}$ such that

$$E \left(\sup_{t \in [0, T]} \|M_{n_k}(t) - M_{n_{k+1}}(t)\|^p \right) \leq 4^{-k}.$$

Then for this subsequence,

$$\begin{aligned} & P \left(\sup_{t \in [0, T]} \|M_{n_k}(t) - M_{n_{k+1}}(t)\| > 2^{-k} \right) \\ & \leq 2^k E \left(\sup_{t \in [0, T]} \|M_{n_k}(t) - M_{n_{k+1}}(t)\|^p \right) \leq 2^{-k} \end{aligned}$$

and so, there is a set of measure zero N such that if $\omega \notin N$, then for all k large enough, $\sup_{t \in [0, T]} \|M_{n_k}(t) - M_{n_{k+1}}(t)\| \leq 2^{-k}$ and so for $\omega \notin N$, there exists M continuous such that $M_{n_k}(t) \rightarrow M(t)$ uniformly in $t \in [0, T]$. Thus for each ω , $\|M(t)\|^p \leq \|M_{n_k}(t)\|^p + \varepsilon$ for all $t \in [0, T]$ if k is large enough. Therefore, $\|M(t)\|^p \leq \sup_t \|M_{n_k}(t)\|^p + \varepsilon$ and so $\sup_t \|M(t)\|^p \leq \sup_t \|M_{n_k}(t)\|^p + \varepsilon$ for all t large enough. It follows that for each ω off a set of measure zero, $\sup_t \|M(t)\|^p \leq \liminf_{k \rightarrow \infty} \sup_t \|M_{n_k}(t)\|^p$. Is $M \in \mathcal{M}_T^p$? By Fatou's lemma,

$$\int_{\Omega} \sup_{t \in [0, T]} \|M(t)\|^p dP \leq \liminf_{k \rightarrow \infty} \|M_{n_k}\|_{\mathcal{M}_T^p(E)}^p$$

which is finite because $\{M_n\}$ is a Cauchy sequence. Thus $M \in \mathcal{M}_T^p(E)$. Now also,

$$\begin{aligned} & \left(\int_{\Omega} \sup_t \|M(t) - M_{n_k}(t)\|^p dP \right)^{1/p} \\ & \leq \liminf_{m \rightarrow \infty} \left(\int_{\Omega} \sup_t \|M_{n_m}(t) - M_{n_k}(t)\|^p dP \right)^{1/p} < \varepsilon \end{aligned}$$

if k is large enough because for $m > k$,

$$\begin{aligned} \|M_{n_m} - M_{n_k}\|_{\mathcal{M}_T^p} & \leq \sum_{r=k}^{\infty} \|M_{n_{r+1}} - M_{n_r}\|_{\mathcal{M}_T^p} \\ & \leq \sum_{r=k}^{\infty} (4^{1/p})^{-r} = (4^{1/p})^{-k} / (1 - (4^{1/p})^{-1}). \end{aligned}$$

This shows that every Cauchy sequence has a convergent subsequence and so the original Cauchy sequence also converges. This shows \mathcal{M}_T^p is complete.

It only remains to verify that if each M_n is a martingale, then so is M a martingale. Let $s \leq t$ and let $B \in \mathcal{F}_s$. For each s , $M_n(s) \rightarrow M(s)$ in $L^p(\Omega)$. Then from the above, $\omega \rightarrow M(s)(\omega)$ is \mathcal{F}_s measurable. Then it follows that

$$\begin{aligned} \int_B M(s) dP &= \lim_{n \rightarrow \infty} \int_B M_n(s) dP = \lim_{n \rightarrow \infty} \int_B E(M_n(t) | \mathcal{F}_s) dP \\ &= \lim_{n \rightarrow \infty} \int_B M_n(t) dP = \int_B M(t) dP \end{aligned}$$

and so by definition, $E(M(t) | \mathcal{F}_s) = M(s)$ which shows M is a martingale.

It is clear that if σ is a stopping time, then if $M \in \mathcal{M}_T^p(E)$ so is M^σ and that

$$\|M^\sigma\|_{\mathcal{M}_T^p(E)} \leq \|M\|_{\mathcal{M}_T^p(E)}$$

Thus if $M_n \rightarrow M$ in $\mathcal{M}_T^p(E)$, then $M_n^\sigma \rightarrow M^\sigma$ in $\mathcal{M}_T^p(E)$. ■

Note that if $M_n \rightarrow M$ in $\mathcal{M}_T^p(E)$, this says $\int_{\Omega} \sup_{t \in [0, T]} \|M_n(t) - M(t)\|^p dP = 0$. Hence this would also be true that $\lim_{n \rightarrow \infty} \int_A \sup_{t \in [0, T]} \|M_n(t) - M(t)\|^p dP = 0$ also, whenever A is a measurable set.

Chapter 32

Quadratic Variation

32.1 How to Recognize a Martingale

The main ideas are most easily understood in the special case where it is assumed the martingale is bounded. Then one can extend to more general situations using a localizing sequence of stopping times.

Let $\{M(t)\}$ be a continuous martingale having values in a separable Hilbert space. The idea is to consider the submartingale, $\{\|M(t)\|^2\}$ and write it as the sum of a martingale and an increasing submartingale. An important part of the argument is the following lemma which gives a checkable criterion for a stochastic process to be a martingale.

Lemma 32.1.1 *Let $\{X(t)\}$ be a stochastic process adapted to the filtration $\{\mathcal{F}_t\}$ for $t \geq 0$. Then it is a martingale for the given filtration if for every stopping time σ it follows*

$$E(X(t)) = E(X(\sigma)).$$

In fact, it suffices to check this on stopping times which have two values.

Proof: Let $s < t$ and $A \in \mathcal{F}_s$. Define a stopping time

$$\sigma(\omega) \equiv s \mathcal{X}_A(\omega) + t \mathcal{X}_{A^c}(\omega)$$

This is a stopping time because $[\sigma \leq l] = \Omega \in \mathcal{F}_l$ if $l \geq t$. Also $[\sigma \leq l] = A \in \mathcal{F}_s \subseteq \mathcal{F}_l$ if $l \in [s, t)$ and $[\sigma \leq l] = \emptyset \in \mathcal{F}_l$ if $l < s$. Then by assumption,

$$\begin{aligned} \int_A X(t) dP + \int_{A^c} X(t) dP &= \\ \overbrace{\int X(t) dP}^{\text{by assumption}} &= \int X(\sigma) dP = \int_A X(s) dP + \int_{A^c} X(t) dP \end{aligned}$$

Therefore,

$$\int_A X(t) dP = \int_A X(s) dP$$

and since $X(s)$ is \mathcal{F}_s measurable, it follows $E(X(t)|\mathcal{F}_s) = X(s)$ a.e. and this shows $\{X(t)\}$ is a martingale. ■

Note that if $t \in [0, T]$, it suffices to check the expectation condition for stopping times which have two values no larger than T .

The following lemma will be useful.

Lemma 32.1.2 *Suppose $X_n \rightarrow X$ in $L^1(\Omega, \mathcal{F}, P; E)$ where E is a separable Banach space. Then letting \mathcal{G} be a σ algebra contained in \mathcal{F} ,*

$$E(X_n|\mathcal{G}) \rightarrow E(X|\mathcal{G})$$

in $L^1(\Omega)$.

Proof: This follows from the definitions and Theorem 24.12.1 on Page 702.

$$\begin{aligned} \int_{\Omega} \|E(X|\mathcal{G}) - E(X_n|\mathcal{G})\| dP &= \int_{\Omega} \|E(X_n - X|\mathcal{G})\| dP \\ &\leq \int_{\Omega} E(\|X_n - X\| |\mathcal{G}) dP = \int_{\Omega} \|X_n - X\| dP \blacksquare \end{aligned}$$

The next corollary is like the earlier result which allows you to take a sufficiently measurable function out of the conditional expectation.

Corollary 32.1.3 *Let X, Y be in $L^2(\Omega, \mathcal{F}, P; H)$ where H is a separable Hilbert space and let X be \mathcal{G} measurable where $\mathcal{G} \subseteq \mathcal{F}$. Then*

$$E((X, Y) | \mathcal{G}) = (X, E(Y | \mathcal{G})) \text{ a.e.}$$

Proof: First let $X = a\mathcal{X}_B$ where $B \in \mathcal{G}$. Then for $A \in \mathcal{G}$,

$$\begin{aligned} \int_A E((a\mathcal{X}_B, Y) | \mathcal{G}) dP &= \int_A \mathcal{X}_B E((a, Y) | \mathcal{G}) dP = \int_A \mathcal{X}_B (a, Y) dP \\ &= \int_{A \cap B} (a, Y) dP = \left(a, \int_{A \cap B} Y dP \right) \end{aligned}$$

$$\begin{aligned} \int_A (a\mathcal{X}_B, E(Y | \mathcal{G})) dP &= \int_A \mathcal{X}_B (a, E(Y | \mathcal{G})) dP \\ &= \left(a, \int_A \mathcal{X}_B E(Y | \mathcal{G}) dP \right) = \left(a, \int_{A \cap B} Y dP \right) \end{aligned}$$

It follows that the formula holds for X simple.

Therefore, letting X_n be a sequence of \mathcal{G} measurable simple functions converging pointwise to X and also in $L^2(\Omega)$,

$$E((X_n, Y) | \mathcal{G}) = (X_n, E(Y | \mathcal{G}))$$

Now the desired formula holds from Lemma 32.1.2. \blacksquare

The following is related to something called a martingale transform. It is a lot like what will happen later with the Ito integral.

Maybe it is a good idea to try and give some reason for considering the following. Say you have a bounded variation and adapted function f and you wanted to consider the Stieltjes integral $\int_0^T f dM$. If f is of bounded variation, this Stieltjes integral will exist from the standard theory of Stieltjes integration. In particular, if g is Stieltjes integrable with respect to df then f is Stieltjes integrable with respect to dg and an integration by parts formula holds. Now assuming M is continuous and f is of bounded variation, you would have the existence of $\int_0^T M df$ and so also the existence of $\int_0^T f dM$. Of course you might have different partitions for each different ω . In the following, this is handled by writing a sum of the form

$$\sum_{k \geq 0} (\xi_k, (M(\tau_{k+1} \wedge t) - M(\tau_k \wedge t)))$$

where ξ_k is in \mathcal{F}_{τ_k} and τ_k is a stopping time, the τ_k being an increasing sequence of stopping times having limit ∞ . You could think of this as the value of f at the left end point. Of course what is happening here pertains to Hilbert space, but the inner product is

sufficiently like multiplication to draw the analogy and you would still have an integration by parts formula and the same result on existence of the integral. The theory of Stieltjes integrals is in my single variable advanced calculus book. An early reference to this is Hobson [28]. See also [2]. What is going to happen here is that these Stieltjes sums will end up being a martingale. The following is stated for the more general situation where $M(t)$ is only right continuous.

Proposition 32.1.4 *Let $\{\tau_k\}$ be an increasing sequence of stopping times for the normal filtration $\{\mathcal{F}_t\}$ such that*

$$\lim_{k \rightarrow \infty} \tau_k = \infty, \quad \tau_0 = 0.$$

Also let ξ_k be \mathcal{F}_{τ_k} measurable with values in H , a separable Hilbert space and let $M(t)$ be a right continuous martingale adapted to the normal filtration \mathcal{F}_t which has the property that $M(t) \in L^2(\Omega; H)$ for all t , $M(0) = 0$. Then if $|\xi_k| \leq C$,

$$E \left(\left(\sum_{k \geq 0} (\xi_k, (M(\tau_{k+1} \wedge t) - M(\tau_k \wedge t))) \right)^2 \right) \leq C^2 E \left(\|M(t)\|^2 \right) \quad (32.1)$$

Proof: First of all, the sum converges because eventually $\tau_k \wedge t = t$. Therefore, for large enough k , $M(\tau_{k+1} \wedge t) - M(\tau_k \wedge t) \equiv \Delta M_k = 0$. Consider first the finite sum, $k \leq q$.

$$E \left(\left(\sum_{k=0}^q (\xi_k, \Delta M_k) \right)^2 \right) \quad (32.2)$$

When the sum is multiplied out, you get mixed terms. Consider one of these mixed terms, $j < k$

$$E \left((\xi_k, \Delta M_k) (\xi_j, \Delta M_j) \right)$$

Using Corollary 32.1.3 and Doob's optional sampling theorem, this equals

$$\begin{aligned} E \left(E \left((\xi_k, \Delta M_k) (\xi_j, \Delta M_j) \mid \mathcal{F}_{\tau_k} \right) \right) &= E \left((\xi_j, \Delta M_j) E \left((\xi_k, \Delta M_k) \mid \mathcal{F}_{\tau_k} \right) \right) \\ &= E \left((\xi_j, \Delta M_j) (\xi_k, E(M(\tau_{k+1} \wedge t) - M(\tau_k \wedge t) \mid \mathcal{F}_{\tau_k})) \right) \\ &= E \left((\xi_j, \Delta M_j) (\xi_k, 0) \right) = 0 \end{aligned}$$

Note that in using the optional sampling theorem, the stopping time $\tau_{k+1} \wedge t$ is bounded.

Therefore, the only terms which survive in 32.2 are the non mixed terms and so this expression reduces to

$$\begin{aligned} \sum_{k=0}^q E (\xi_k, \Delta M_k)^2 &\leq C^2 \sum_{k=0}^q E \left(\|\Delta M_k\|^2 \right) \\ &= C^2 \sum_{k=0}^q E \left(\|M(\tau_{k+1} \wedge t) - M(\tau_k \wedge t)\|^2 \right) \end{aligned}$$

$$\begin{aligned}
&= C^2 \sum_{k=0}^q E \left(\|M(\tau_{k+1} \wedge t)\|^2 \right) + E \left(\|M(\tau_k \wedge t)\|^2 \right) \\
&\quad - 2E \left((M(\tau_k \wedge t), M(\tau_{k+1} \wedge t)) \right)
\end{aligned} \tag{32.3}$$

Consider the term $E \left((M(\tau_k \wedge t), M(\tau_{k+1} \wedge t)) \right)$. By Doob's optional sampling theorem for martingales and Corollary 32.1.3 again, this equals

$$\begin{aligned}
&E \left(E \left((M(\tau_k \wedge t), M(\tau_{k+1} \wedge t)) \mid \mathcal{F}_{\tau_k} \right) \right) \\
&= E \left((M(\tau_k \wedge t), E(M(\tau_{k+1} \wedge t) \mid \mathcal{F}_{\tau_k})) \right) \\
&= E \left((M(\tau_k \wedge t), M(\tau_{k+1} \wedge t \wedge \tau_k)) \right) = E \left(\|M(\tau_k \wedge t)\|^2 \right)
\end{aligned}$$

It follows 32.3 equals

$$\begin{aligned}
&C^2 \sum_{k=0}^q E \left(\|M(\tau_{k+1} \wedge t)\|^2 \right) - E \left(\|M(\tau_k \wedge t)\|^2 \right) \\
&\leq C^2 E \left(\|M(t)\|^2 \right).
\end{aligned}$$

Then from Fatou's lemma,

$$\begin{aligned}
&E \left(\left(\sum_{k=0}^q (\xi_k, (M(\tau_{k+1} \wedge t) - M(\tau_k \wedge t))) \right)^2 \right) \leq \\
&\liminf_{q \rightarrow \infty} E \left(\left(\sum_{k=0}^q (\xi_k, (M(\tau_{k+1} \wedge t) - M(\tau_k \wedge t))) \right)^2 \right) \\
&\leq C^2 E \left(\|M(t)\|^2 \right) \blacksquare
\end{aligned}$$

Now here is an interesting lemma which will be used to prove uniqueness in the main result.

32.2 Martingales and Total Variation

Lemma 32.2.1 *Let \mathcal{F}_t be a normal filtration and let $A(t), B(t)$ be adapted to \mathcal{F}_t , continuous, and increasing with $A(0) = B(0) = 0$ and suppose $A(t) - B(t)$ is a martingale. Then $A(t) - B(t) = 0$ for all t .*

Proof: I shall show $A(l) = B(l)$ where l is arbitrary. Let $M(t)$ be the name of the martingale. Define a stopping time

$$\begin{aligned}
\tau &\equiv \inf \{t > 0 : |M(t)| > C\} \wedge l \wedge \inf \{t > 0 : A(t) > C\} \\
&\quad \wedge \inf \{t > 0 : B(t) > C\}
\end{aligned}$$

where $\inf(\emptyset) \equiv \infty$ and denote the stopped martingale $M^\tau(t) \equiv M(t \wedge \tau)$. Then this is also a martingale with respect to the filtration \mathcal{F}_t because by Doob's optional sampling theorem for martingales. Recall why this is: if $s < t$,

$$E(M^\tau(t) \mid \mathcal{F}_s) \equiv E(M(\tau \wedge t) \mid \mathcal{F}_s) = M(\tau \wedge t \wedge s) = M(\tau \wedge s) = M^\tau(s)$$

Note the bounded stopping time is $\tau \wedge t$ and the other one is $\sigma = s$ in this theorem. Then M^τ is a continuous martingale which is also uniformly bounded. It equals $A^\tau - B^\tau$. The stopping time ensures A^τ and B^τ are uniformly bounded by C . Thus all of $|M^\tau(t)|, B^\tau(t), A^\tau(t)$ are bounded by C on $[0, l]$. Now let $\mathcal{P}_n = \{t_k\}_{k=1}^n$ be a uniform partition of $[0, l]$ and let $M^\tau(\mathcal{P}_n)$ denote

$$M^\tau(\mathcal{P}_n) \equiv \max \{|M^\tau(t_{i+1}) - M^\tau(t_i)|\}_{i=1}^n.$$

Then

$$E(M^\tau(l)^2) = E\left(\left(\sum_{k=0}^{n-1} M^\tau(t_{k+1}) - M^\tau(t_k)\right)^2\right)$$

Now consider a mixed term in the sum where $j < k$.

$$\begin{aligned} & E((M^\tau(t_{k+1}) - M^\tau(t_k))(M^\tau(t_{j+1}) - M^\tau(t_j))) \\ &= E(E((M^\tau(t_{k+1}) - M^\tau(t_k))(M^\tau(t_{j+1}) - M^\tau(t_j)) | \mathcal{F}_{t_k})) \\ &= E((M^\tau(t_{j+1}) - M^\tau(t_j)) E((M^\tau(t_{k+1}) - M^\tau(t_k)) | \mathcal{F}_{t_k})) \\ &= E((M^\tau(t_{j+1}) - M^\tau(t_j))(M^\tau(t_k) - M^\tau(t_k))) = 0 \end{aligned}$$

It follows

$$\begin{aligned} E(M^\tau(l)^2) &= E\left(\sum_{k=0}^{n-1} (M^\tau(t_{k+1}) - M^\tau(t_k))^2\right) \\ &\leq E\left(\sum_{k=0}^{n-1} M^\tau(\mathcal{P}_n) |M^\tau(t_{k+1}) - M^\tau(t_k)|\right) \\ &\leq E\left(\sum_{k=0}^{n-1} M^\tau(\mathcal{P}_n) (|A^\tau(t_{k+1}) - A^\tau(t_k)| + |B^\tau(t_{k+1}) - B^\tau(t_k)|)\right) \\ &\leq E\left(M^\tau(\mathcal{P}_n) \sum_{k=0}^{n-1} (|A^\tau(t_{k+1}) - A^\tau(t_k)| + |B^\tau(t_{k+1}) - B^\tau(t_k)|)\right) \\ &\leq E(M^\tau(\mathcal{P}_n) 2C) \end{aligned}$$

the last step holding because A and B are increasing. Now letting $n \rightarrow \infty$, the right side converges to 0 by the dominated convergence theorem and $\lim_{n \rightarrow \infty} M^\tau(\mathcal{P}_n)(\omega) = 0$ because of continuity of M . Thus for $\tau = \tau_C$ given above, $M(l \wedge \tau_C) = 0$ a.e. Now let $C \in \mathbb{N}$ and let N_C be the exceptional set off which $M(l \wedge \tau_C) = 0$. Then letting N_l denote the union of all these exceptional sets for $C \in \mathbb{N}$, it is also a set of measure zero and for ω not in this set, $M(l \wedge \tau_C) = 0$ for all C . Since the martingale is continuous, it follows for each such ω , eventually $\tau_C > l$ and so $M(l) = 0$. Thus for $\omega \notin N_l, M(l)(\omega) = 0$. Now let $N = \cup_{l \in \mathbb{Q} \cap [0, \infty)} N_l$. Then for $\omega \notin N, M(l)(\omega) = 0$ for all $l \in \mathbb{Q} \cap [0, \infty)$ and so by continuity, this is true for all positive l . ■

Note this shows a continuous martingale is not of bounded variation unless it is a constant.

If you had a continuous bounded variation function $f(t)$, you might want to do something like $\int_0^T \frac{1}{2} (f'(t), f(t)) dt = \int_0^T \frac{1}{2} (f(t), df) = |f(T)|^2$. We do this all the time when

we discuss curves. From what was just shown, however, it will not be possible to do this directly in the context of Stieltjes integrals. What happens is something else. We get something called the quadratic variation of the martingale M denoted by $[M]$ which is increasing and $\|M\|^2 = [M] + N$ where N is a martingale. It will make perfect sense to write $\int_0^t f(s) d[M](s)$. Thus N does not have finite total variation but $[M]$ does.

32.3 The Quadratic Variation

This section is on the quadratic variation of a martingale. Actually, you can also consider the quadratic variation of a local martingale which is more general. Therefore, this concept is defined first. We will generally assume $M(0) = 0$ since there is no real loss of generality in doing so. One can simply subtract $M(0)$ otherwise. What is about to be presented is called the quadratic variation because it considers the variation of $\|M(t)\|^2$ rather than $M(t)$ which, as just shown is not finite.

Definition 32.3.1 *Let $\{M(t)\}$ be adapted to the normal filtration \mathcal{F}_t for $t > 0$. Then $\{M(t)\}$ is a local martingale (submartingale) if there exist stopping times τ_n increasing to infinity such that for each n , the process $M^{\tau_n}(t) \equiv M(t \wedge \tau_n)$ is a martingale (submartingale) with respect to the given filtration. The sequence of stopping times is called a localizing sequence. The martingale M^{τ_n} is called the stopped martingale. Exactly the same convention applies to a localized submartingale. When M is continuous, we can always assume that M^{τ_n} is a bounded martingale (submartingale) by taking the minimum of τ_n with the first hitting time of $n \in \mathbb{N}$ by $\|M\|$. I will use this observation whenever convenient. If this is done, then $\tau_n = \infty$ for n large enough.*

Observation 32.3.2 *If M is a local martingale (submartingale) and if σ is a stopping time, then M^σ is also a local martingale (submartingale).*

To see this, use the localizing sequence. Say M is a local submartingale. The case of a martingale is similar. Let $s \leq t$.

$$E((M^\sigma)^{\tau_n}(t) | \mathcal{F}_s) = E(M^{\tau_n}(\sigma \wedge t) | \mathcal{F}_s) \geq M^{\tau_n}(\sigma \wedge s) = (M^\sigma)^{\tau_n}(s).$$

By the optional sampling theorem.

Proposition 32.3.3 *If $M(t)$ is a continuous local martingale (submartingale) for a normal filtration as above, $M(0) = 0$, then there exists a localizing sequence τ_n such that for each n the stopped martingale (submartingale) M^{τ_n} is uniformly bounded. Also if M is a martingale, then M^τ is also a martingale (submartingale). If τ_n is an increasing sequence of stopping times such that $\lim_{n \rightarrow \infty} \tau_n = \infty$, and for each τ_n and real valued stopping time δ , there exists a function X of $\tau_n \wedge \delta$ such that $X(\tau_n \wedge \delta)$ is $\mathcal{F}_{\tau_n \wedge \delta}$ measurable, then $\lim_{n \rightarrow \infty} X(\tau_n \wedge \delta) \equiv X(\delta)$ exists for each ω and $X(\delta)$ is \mathcal{F}_δ measurable.*

Proof: First review the claim about M^τ being a martingale (submartingale) when M is. By optional sampling theorem,

$$E(M^\tau(t) | \mathcal{F}_s) = E(M(\tau \wedge t) | \mathcal{F}_s) = M(\tau \wedge t \wedge s) = M^\tau(s).$$

The case where M is a submartingale is similar.

Next suppose σ_n is a localizing sequence for the local martingale(submartingale) M . Then define

$$\eta_n \equiv \inf\{t > 0 : \|M(t)\| > n\}.$$

Therefore, by continuity of M , $\|M(\eta_n)\| \leq n$. Now consider $\tau_n \equiv \eta_n \wedge \sigma_n$. This is an increasing sequence of stopping times. By continuity of M , it must be the case that $\eta_n \rightarrow \infty$. Hence $\sigma_n \wedge \eta_n \rightarrow \infty$.

Finally, consider the last claim. Pick ω . Then $X(\tau_n(\omega) \wedge \delta(\omega))(\omega)$ is eventually constant as $n \rightarrow \infty$ because for all n large enough, $\tau_n(\omega) > \delta(\omega)$ and so this sequence of functions converges pointwise. That which it converges to, denoted by $X(\delta)$, is \mathcal{F}_δ measurable because each function $\omega \rightarrow X(\tau_n(\omega) \wedge \delta(\omega))(\omega)$ is $\mathcal{F}_{\delta \wedge \tau_n} \subseteq \mathcal{F}_\delta$ measurable. ■

Observation 32.3.4 Suppose M is a local martingale and τ_n is a localizing sequence of stoppings times. Does $M^{\tau_n}(t)$ converge in probability to $M(t)$? $M^{\tau_k}(t) = M(t)$ at ω where $\tau_k(\omega) = \infty$ and so $\{ \|M^{\tau_k}(t) - M(t)\| > \varepsilon \} \subseteq [\tau_k < \infty]$ and $P([\tau_k < \infty]) \rightarrow 0$ by assumption that τ_k is a localizing sequence.

One can also give a generalization of Lemma 32.2.1 to conclude a local martingale must be constant or else they must fail to be of bounded variation.

Corollary 32.3.5 Let \mathcal{F}_t be a normal filtration and let $A(t), B(t)$ be adapted to \mathcal{F}_t , continuous, and increasing with $A(0) = B(0) = 0$ and suppose $A(t) - B(t) \equiv M(t)$ is a local martingale. Then $M(t) = A(t) - B(t) = 0$ a.e. for all t .

Proof: Let $\{\tau_n\}$ be a localizing sequence for M . For given n , consider the martingale,

$$M^{\tau_n}(t) = A^{\tau_n}(t) - B^{\tau_n}(t)$$

Then from Lemma 32.2.1, it follows $M^{\tau_n}(t) = 0$ for all t for all $\omega \notin N_n$, a set of measure 0. Let $N = \cup_n N_n$. Then for $\omega \notin N$, $M(\tau_n(\omega) \wedge t)(\omega) = 0$. Let $n \rightarrow \infty$ to conclude that $M(t)(\omega) = 0$. Therefore, $M(t)(\omega) = 0$ for all t . ■

Recall Example 31.3.13 on Page 841. For convenience, here is a version of what it says.

Lemma 32.3.6 Let $X(t)$ be continuous and adapted to a normal filtration \mathcal{F}_t and let η be a stopping time. Then if K is a closed set,

$$\tau \equiv \inf\{t > \eta : X(t) \in K\}$$

is also a stopping time.

Proof: First consider $Y(t) = X(t \vee \eta) - X(\eta)$. I claim that $Y(t)$ is adapted to \mathcal{F}_t . Consider U and open set and $[Y(t) \in U]$. Is it in \mathcal{F}_t ? We know it is in $\mathcal{F}_{t \vee \eta}$. It equals

$$([Y(t) \in U] \cap [\eta \leq t]) \cup ([Y(t) \in U] \cap [\eta > t])$$

Consider the second of these sets. It equals

$$([X(\eta) - X(\eta) \in U] \cap [\eta > t])$$

If $0 \in U$, then it reduces to $[\eta > t] \in \mathcal{F}_t$. If $0 \notin U$, then it reduces to \emptyset still in \mathcal{F}_t . Next consider the first set. It equals

$$\begin{aligned} & [X(t \vee \eta) - X(\eta) \in U] \cap [\eta \leq t] \\ = & [X(t \vee \eta) - X(\eta) \in U] \cap [t \vee \eta \leq t] \in \mathcal{F}_t \end{aligned}$$

from the definition of $\mathcal{F}_{t \vee \eta}$. (You know that $[X(t \vee \eta) - X(\eta) \in U] \in \mathcal{F}_{t \vee \eta}$ and so when this is intersected with $[t \vee \eta \leq t]$ one obtains a set in \mathcal{F}_t . This is what it means to be in $\mathcal{F}_{t \vee \eta}$.) Now τ is just the first hitting time of $Y(t)$ of the closed set. ■

Proposition 32.3.7 *Let $M(t)$ be a continuous local martingale for $t \in [0, T]$ having values in H a separable Hilbert space adapted to the normal filtration $\{\mathcal{F}_t\}$ such that $M(0) = 0$. Then there exists a unique continuous, increasing, nonnegative, local submartingale $[M](t)$ called the quadratic variation such that*

$$\|M(t)\|^2 - [M](t)$$

is a real local martingale and $[M](0) = 0$. Here $t \in [0, T]$. If δ is any stopping time

$$[M^\delta] = [M]^\delta$$

Proof: First it is necessary to define some stopping times. Define stopping times $\tau_0^n \equiv \eta_0^n \equiv 0$.

$$\begin{aligned} \eta_{k+1}^n & \equiv \inf \{s > \eta_k^n : \|M(s) - M(\eta_k^n)\| = 2^{-n}\}, \\ \tau_k^n & \equiv \eta_k^n \wedge T \end{aligned}$$

where $\inf \emptyset \equiv \infty$. These are stopping times by Example 31.3.13 on Page 841. See also the above Lemma 32.3.6. Then for $t > 0$ and δ any stopping time, and fixed ω , for some k ,

$$t \wedge \delta \in I_k(\omega), I_0(\omega) \equiv [\tau_0^n(\omega), \tau_1^n(\omega)], I_k(\omega) \equiv (\tau_k^n(\omega), \tau_{k+1}^n(\omega)] \text{ some } k$$

Here is why. The sequence $\{\tau_k^n(\omega)\}_{k=1}^\infty$ eventually equals T for all n sufficiently large. This is because if it did not, it would converge, being bounded above by T and then by continuity of M , $\{M(\tau_k^n(\omega))\}_{k=1}^\infty$ would be a Cauchy sequence contrary to the requirement that

$$\begin{aligned} & \|M(\tau_{k+1}^n(\omega)) - M(\tau_k^n(\omega))\| \\ = & \|M(\eta_{k+1}^n(\omega)) - M(\eta_k^n(\omega))\| = 2^{-n}. \end{aligned}$$

Note that if δ is any stopping time, then

$$\begin{aligned} & \|M(t \wedge \delta \wedge \tau_{k+1}^n) - M(t \wedge \delta \wedge \tau_k^n)\| \\ = & \|M^\delta(t \wedge \tau_{k+1}^n) - M^\delta(t \wedge \tau_k^n)\| \leq 2^{-n} \end{aligned}$$

You can see this is the case by considering the cases, $t \wedge \delta \geq \tau_{k+1}^n$, $t \wedge \delta \in [\tau_k^n, \tau_{k+1}^n)$, and $t \wedge \delta < \tau_k^n$. It is only this approximation property and the fact that the τ_k^n partition $[0, T]$ which is important in the following argument.

Now let α_n be a localizing sequence such that M^{α_n} is bounded as in Proposition 32.3.3. Thus $M^{\alpha_n}(t) \in L^2(\Omega)$ and this is all that is needed. In what follows, let δ be a stopping

time and denote $M^{\alpha_p \wedge \delta}$ by M to save notation. Thus M will be uniformly bounded and from the definition of the stopping times τ_k^n , for $t \in [0, T]$,

$$M(t) \equiv \sum_{k \geq 0} M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n), \quad (32.4)$$

and the terms of the series are eventually 0, as soon as $\tau_k^n = \infty$.

Therefore,

$$\|M(t)\|^2 = \left\| \sum_{k \geq 0} M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n) \right\|^2$$

Then this equals

$$\begin{aligned} &= \sum_{k \geq 0} \|M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n)\|^2 \\ &+ \sum_{j \neq k} ((M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n)), (M(t \wedge \tau_{j+1}^n) - M(t \wedge \tau_j^n))) \end{aligned} \quad (32.5)$$

Consider the second sum. It equals

$$\begin{aligned} &2 \sum_{k \geq 0} \sum_{j=0}^{k-1} ((M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n)), (M(t \wedge \tau_{j+1}^n) - M(t \wedge \tau_j^n))) \\ &= 2 \sum_{k \geq 0} \left((M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n)), \overset{\text{telescopes}}{\sum_{j=0}^{k-1} (M(t \wedge \tau_{j+1}^n) - M(t \wedge \tau_j^n))} \right) \\ &= 2 \sum_{k \geq 0} ((M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n)), M(t \wedge \tau_k^n)) \end{aligned}$$

This last sum equals $P_n(t)$ defined as

$$2 \sum_{k \geq 0} (M(\tau_k^n), (M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n))) \equiv P_n(t) \quad (32.6)$$

This is because in the k^{th} term, if $t \geq \tau_k^n$, then it reduces to

$$(M(\tau_k^n), (M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n)))$$

while if $t < \tau_k^n$, then the term reduces to $((M(t) - M(t)), M(t)) = 0$ which is also the same as

$$(M(\tau_k^n), (M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n))).$$

This is a finite sum because eventually, for large enough k , $\tau_k^n = T$. However the number of nonzero terms depends on ω . This is not a good thing. However, a little more can be said. In fact the sum in 32.6 converges in $L^2(\Omega)$. Say $\|M(t, \omega)\| \leq C$.

$$E \left(\left(\sum_{k \geq p}^q (M(\tau_k^n), (M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n))) \right)^2 \right)$$

$$= \sum_{k \geq p}^q E \left(\left(M(\tau_k^n), (M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n)) \right)^2 \right) + \text{mixed terms} \quad (32.7)$$

Consider one of these mixed terms for $j < k$.

$$E \left(\left(M(\tau_j^n), \overbrace{M(t \wedge \tau_{j+1}^n) - M(t \wedge \tau_j^n)}^{\Delta_j} \right) \left(M(\tau_k^n), \overbrace{M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n)}^{\Delta_k} \right) \right)$$

Then it equals

$$\begin{aligned} & E \left(E \left((M(\tau_j^n), \Delta_j) (M(\tau_k^n), \Delta_k) \mid \mathcal{F}_{\tau_k} \right) \right) \\ &= E \left((M(\tau_j^n), \Delta_j) E \left((M(\tau_k^n), \Delta_k) \mid \mathcal{F}_{\tau_k} \right) \right) \\ &= E \left((M(\tau_j^n), \Delta_j) (M(\tau_k^n), E(\Delta_k \mid \mathcal{F}_{\tau_k})) \right) = 0 \end{aligned}$$

since $E(\Delta_k \mid \mathcal{F}_{\tau_k}) = E(M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n) \mid \mathcal{F}_{\tau_k}) = 0$. Now since the mixed terms equal 0, it follows from 32.7, that expression is dominated by

$$C^2 \sum_{k \geq p}^q E \left(\|M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n)\|^2 \right) \quad (32.8)$$

A mixed term in the above is of the form: For $j < k$,

$$E(\Delta_k, \Delta_j) = E \left(E \left((\Delta_k, \Delta_j) \mid \mathcal{F}_{\tau_k} \right) \right) = E \left((\Delta_j, E(\Delta_k \mid \mathcal{F}_{\tau_k})) \right) = 0$$

Thus 32.8 equals

$$\begin{aligned} & C^2 \sum_{k=p}^q E \left(\|M(t \wedge \tau_{k+1}^n)\|^2 \right) - E \left(\|M(t \wedge \tau_k^n)\|^2 \right) \\ &= C^2 E \left(\|M(t \wedge \tau_{q+1}^n)\|^2 - \|M(t \wedge \tau_p^n)\|^2 \right) \end{aligned}$$

The integrand converges to 0 as $p, q \rightarrow \infty$ and the uniform bound on M allows a use of the dominated convergence theorem. Thus the partial sums of the series of 32.6 converge in $L^2(\Omega)$ as claimed.

By adding in the values of $\{\tau_k^{n+1}\}$ $P_n(t)$ can be written in the form

$$2 \sum_{k \geq 0} (M(\tau_k^{n+1}), (M(t \wedge \tau_{k+1}^{n+1}) - M(t \wedge \tau_k^{n+1})))$$

where $\tau_k^{n+1'}$ has some repeats. From the construction,

$$\|M(\tau_k^{n+1'}) - M(\tau_k^{n+1})\| \leq 2^{-(n+1)}$$

Thus

$$P_n(t) - P_{n+1}(t) = 2 \sum_{k \geq 0} (M(\tau_k^{n+1'}) - M(\tau_k^{n+1}), (M(t \wedge \tau_{k+1}^{n+1}) - M(t \wedge \tau_k^{n+1})))$$

and so from Proposition 32.1.4 applied to $\xi_k \equiv M(\tau_k^{n+1}) - M(\tau_k^n)$,

$$E\left(\|P_n(t) - P_{n+1}(t)\|^2\right) \leq \left(2^{-2n} E\left(\|M(t)\|^2\right)\right). \quad (32.9)$$

Now $t \rightarrow P_n(t)$ is continuous because it is a finite sum of continuous functions. It is also the case that $\{P_n(t)\}$ is a martingale. To see this use Lemma 32.1.1. Let σ be a stopping time having two values. Then using Corollary 32.1.3 and the Doob optional sampling theorem, Theorem 31.3.16

$$\begin{aligned} & E\left(\sum_{k=0}^q (M(\tau_k^n), (M(\sigma \wedge \tau_{k+1}^n) - M(\sigma \wedge \tau_k^n)))\right) \\ &= \sum_{k=0}^q E\left((M(\tau_k^n), (M(\sigma \wedge \tau_{k+1}^n) - M(\sigma \wedge \tau_k^n)))\right) \\ &= \sum_{k=0}^q E\left(\left(E(M(\tau_k^n), (M(\sigma \wedge \tau_{k+1}^n) - M(\sigma \wedge \tau_k^n)) | \mathcal{F}_{\tau_k^n})\right)\right) \\ &= \sum_{k=0}^q E\left(\left(M(\tau_k^n), E(M(\sigma \wedge \tau_{k+1}^n) - M(\sigma \wedge \tau_k^n) | \mathcal{F}_{\tau_k^n})\right)\right) \\ &= \sum_{k=0}^q E\left((M(\tau_k^n), E(M(\sigma \wedge \tau_{k+1}^n \wedge \tau_k^n) - M(\sigma \wedge \tau_k^n)))\right) = 0 \end{aligned}$$

Note the Doob theorem applies because $\sigma \wedge \tau_{k+1}^n$ is a bounded stopping time due to the fact σ has only two values. Similarly

$$\begin{aligned} & E\left(\sum_{k=0}^q (M(\tau_k^n), (M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n)))\right) \\ &= \sum_{k=0}^q E\left((M(\tau_k^n), (M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n)))\right) \\ &= \sum_{k=0}^q E\left(\left(E(M(\tau_k^n), (M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n)) | \mathcal{F}_{\tau_k^n})\right)\right) \\ &= \sum_{k=0}^q E\left(\left(M(\tau_k^n), E(M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n) | \mathcal{F}_{\tau_k^n})\right)\right) \\ &= \sum_{k=0}^q E\left((M(\tau_k^n), E(M(t \wedge \tau_{k+1}^n \wedge \tau_k^n) - M(t \wedge \tau_k^n)))\right) = 0 \end{aligned}$$

It follows each partial sum for $P_n(t)$ is a martingale. As shown above, these partial sums converge in $L^2(\Omega)$ and so it follows that $P_n(t)$ is also a martingale. Note the Doob theorem applies because $t \wedge \tau_{k+1}^n$ is a bounded stopping time.

I want to argue that P_n is a Cauchy sequence in $\mathcal{M}_T^2(\mathbb{R})$. By Theorem 31.4.3 and continuity of P_n which yields appropriate measurability in $\sup_{t \leq T} |P_n(t) - P_{n+1}(t)|$,

$$E\left(\left(\sup_{t \leq T} |P_n(t) - P_{n+1}(t)|\right)^2\right)^{1/2} \leq 2E\left(|P_n(T) - P_{n+1}(T)|^2\right)^{1/2}$$

By 32.9, $\leq 2^{-n} E \left(\|M(T)\|^2 \right)^{1/2}$ which shows $\{P_n\}$ is a Cauchy sequence in $\mathcal{M}_T^2(\mathbb{R})$.

Therefore, by Proposition 31.7.2, there exists $\{N(t)\} \in \mathcal{M}_T^2(\mathbb{R})$ such that $P_n \rightarrow N$ in $\mathcal{M}_T^2(H)$. That is

$$\lim_{n \rightarrow \infty} E \left(\sup_{t \in [0, T]} |P_n(t) - N(t)|^2 \right)^{1/2} = 0.$$

Since $\{N(t)\} \in \mathcal{M}_T^2(\mathbb{R})$, it is a continuous martingale and $N(t) \in L^2(\Omega)$, and $N(0) = 0$ because this is true of each $P_n(0)$. From the above 32.5,

$$\|M(t)\|^2 = Q_n(t) + P_n(t) \quad (32.10)$$

where

$$Q_n(t) = \sum_{k \geq 0} \|M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n)\|^2$$

and $P_n(t)$ is a martingale. Then from 32.10, $Q_n(t)$ is a submartingale and converges for each t to something, denoted as $[M](t)$ in $L^1(\Omega)$ uniformly in $t \in [0, T]$. This is because $P_n(t)$ converges uniformly on $[0, T]$ to $N(t)$ in $L^2(\Omega)$ and $\|M(t)\|^2$ does not depend on n . Then also $[M]$ is a submartingale which equals 0 at 0 because this is true of Q_n and because if $A \in \mathcal{F}_s$ where $s < t$,

$$\begin{aligned} \int_A E([M](t) | \mathcal{F}_s) dP &\equiv \int_A [M](t) dP = \lim_{n \rightarrow \infty} \int_A (\|M(t)\|^2 - P_n(t)) dP \\ &= \lim_{n \rightarrow \infty} \int_A E(\|M(t)\|^2 - P_n(t) | \mathcal{F}_s) dP \geq \liminf_{n \rightarrow \infty} \int_A \|M(s)\|^2 - P_n(s) dP \\ &= \liminf_{n \rightarrow \infty} \int_A Q_n(s) dP \geq \int_A [M](s) dP. \end{aligned}$$

Note that $Q_n(t)$ is increasing because as t increases, the definition allows for the possibility of more nonzero terms in the sum. Therefore, $[M](t)$ is also increasing in t . The function $t \rightarrow [M](t)$ is continuous because $\|M(t)\|^2 = [M](t) + N(t)$ and $t \rightarrow N(t)$ is continuous as is $t \rightarrow \|M(t)\|^2$. That is, off a set of measure zero, these are both continuous functions of t and so the same is true of $[M]$.

Now put back in $M^{\alpha_p \wedge \delta}$ in place of M where δ is a stopping time. From the above, this has shown

$$\|M^{\alpha_p \wedge \delta}(t)\|^2 = [M^{\alpha_p \wedge \delta}](t) + N_p(t)$$

where N_p is a martingale and

$$\begin{aligned} [M^{\alpha_p \wedge \delta}](t) &= \lim_{n \rightarrow \infty} \sum_{k \geq 0} \|M^{\alpha_p \wedge \delta}(t \wedge \tau_{k+1}^n) - M^{\alpha_p \wedge \delta}(t \wedge \tau_k^n)\|^2 \\ &= \lim_{n \rightarrow \infty} \sum_{k \geq 0} \|M(t \wedge \tau_{k+1}^n \wedge \alpha_p \wedge \delta) - M(t \wedge \tau_k^n \wedge \alpha_p \wedge \delta)\|^2 \text{ in } L^1(\Omega), \end{aligned} \quad (32.11)$$

the convergence being uniform on $[0, T]$. The above formula shows that $[M^{\alpha_p \wedge \delta}](t)$ is a $\mathcal{F}_{t \wedge \delta \wedge \alpha_p}$ measurable random variable which depends on $t \wedge \delta \wedge \alpha_p$. (Note that $t \wedge \delta$ is a real valued stopping time even if $\delta = \infty$.) Therefore, by Proposition 32.3.3, there exists a random variable, denoted as $[M^\delta](t)$ which is the pointwise limit as $p \rightarrow \infty$ of these random

variables which is $\mathcal{F}_{t \wedge \delta}$ measurable because, for a given ω , when α_p becomes larger than t , the sum in 32.11 loses its dependence on p . Thus from pointwise convergence in 32.11,

$$[M^\delta](t) \equiv \lim_{n \rightarrow \infty} \sum_{k \geq 0} \|M(t \wedge \delta \wedge \tau_{k+1}^n) - M(t \wedge \delta \wedge \tau_k^n)\|^2$$

In case $\delta = \infty$, the above gives an \mathcal{F}_t measurable random variable denoted by $[M](t)$ such that

$$[M](t) \equiv \lim_{n \rightarrow \infty} \sum_{k \geq 0} \|M(t \wedge \tau_{k+1}^n) - M(t \wedge \tau_k^n)\|^2$$

Now stopping with the stopping time δ , this shows that

$$[M^\delta](t) \equiv \lim_{n \rightarrow \infty} \sum_{k \geq 0} \|M(t \wedge \delta \wedge \tau_{k+1}^n) - M(t \wedge \delta \wedge \tau_k^n)\|^2 = [M]^\delta(t)$$

That is, the quadratic variation of the stopped local martingale makes sense a.e. and equals the stopped quadratic variation of the local martingale.

This has now shown that

$$\begin{aligned} \|M^{\alpha_n}(t)\|^2 - [M]^{\alpha_n}(t) &= \|M^{\alpha_n}(t)\|^2 - [M^{\alpha_n}](t) \\ &= N_n(t), \quad N_n(t) \text{ a martingale} \end{aligned}$$

and both of the random variables on the left converge pointwise as $n \rightarrow \infty$ to a function which is \mathcal{F}_t measurable. Hence so does $N_n(t)$. Of course $N_n(t)$ is likewise a function of $\alpha_n \wedge t$ and so by Proposition 32.3.3 again, it converges pointwise to a \mathcal{F}_t measurable function called $N(t)$ and $N(t)$ is a continuous local martingale.

It remains to consider the claim about the uniqueness. Suppose then there are two which work, $[M]$, and $[M]_1$. Then $[M] - [M]_1$ equals a local martingale G which is 0 when $t = 0$. Thus the uniqueness assertion follows from Corollary 32.3.5. ■

Here is a corollary which tells how to manipulate stopping times. It is contained in the above proposition, but it is worth emphasizing it from a different point of view.

Corollary 32.3.8 *In the situation of Proposition 32.3.7 let τ be a stopping time. Then*

$$[M^\tau] = [M]^\tau.$$

Proof:

$$[M]^\tau(t) + N_1(t) = \left(\|M\|^2 \right)^\tau(t) = \|M^\tau\|^2(t) = [M^\tau](t) + N_2(t)$$

where N_i is a local martingale. Therefore,

$$[M]^\tau(t) - [M^\tau](t) = N_2(t) - N_1(t),$$

a local martingale. Therefore, by Corollary 32.3.5, this shows $[M]^\tau(t) - [M^\tau](t) = 0$. ■

32.4 The Covariation

Definition 32.4.1 *The covariation of two continuous H valued local martingales for H a separable Hilbert space $M, N, M(0) = 0 = N(0)$, is defined as follows.*

$$[M, N] \equiv \frac{1}{4} ([M + N] - [M - N])$$

Lemma 32.4.2 *The following hold for the covariation.*

$$[M] = [M, M]$$

$$\begin{aligned} [M, N] &= \text{local martingale} + \frac{1}{4} \left(\|M + N\|^2 - \|M - N\|^2 \right) \\ &= (M, N) + \text{local martingale}. \end{aligned}$$

Proof: From the definition of covariation,

$$\begin{aligned} [M] &= \|M\|^2 - \mathcal{N}_1 \\ [M, M] &= \frac{1}{4} ([M + M] - [M - M]) = \frac{1}{4} \left(\|M + M\|^2 - \mathcal{N}_2 \right) \\ &= \|M\|^2 - \frac{1}{4} \mathcal{N}_2 \end{aligned}$$

where \mathcal{N}_i is a local martingale. Thus $[M] - [M, M]$ is equal to the difference of two increasing continuous adapted processes and it also equals a local martingale. By Corollary 32.3.5, this process must equal 0. Now consider the second claim.

$$\begin{aligned} [M, N] &= \frac{1}{4} ([M + N] - [M - N]) = \frac{1}{4} \left(\|M + N\|^2 - \|M - N\|^2 + \mathcal{N} \right) \\ &= (M, N) + \frac{1}{4} \mathcal{N} \end{aligned}$$

where \mathcal{N} is a local martingale. ■

Corollary 32.4.3 *Let M, N be two continuous local martingales,*

$$M(0) = N(0) = 0,$$

as in Proposition 32.3.7. Then $[M, N]$ is of bounded variation and

$$(M, N)_H - [M, N]$$

is a local martingale. Also for τ a stopping time,

$$[M, N]^\tau = [M^\tau, N^\tau] = [M^\tau, N] = [M, N^\tau].$$

In addition to this,

$$[M - M^\tau] = [M] - [M^\tau] \leq [M]$$

and also

$$M, N \rightarrow [M, N]$$

is bilinear and symmetric.

Proof: Since $[M, N]$ is the difference of increasing functions, it is of bounded variation.

$$\begin{aligned} (M, N)_H - [M, N] &= \overbrace{\frac{1}{4} \left(\|M + N\|^2 - \|M - N\|^2 \right)}^{(M, N)_H} \\ &\quad - \overbrace{\frac{1}{4} ([M + N] - [M - N])}^{[M, N]} \end{aligned}$$

which equals a local martingale from the definition of $[M + N]$ and $[M - N]$. It remains to verify the claim about the stopping time. Using Corollary 32.3.8

$$\begin{aligned} [M, N]^\tau &= \frac{1}{4} ([M + N] - [M - N])^\tau \\ &= \frac{1}{4} ([M + N]^\tau - [M - N]^\tau) \\ &= \frac{1}{4} ([M^\tau + N^\tau] - [M^\tau - N^\tau]) \equiv [M^\tau, N^\tau]. \end{aligned}$$

The really interesting part is the next equality. This will involve Corollary 32.3.5.

$$\begin{aligned} [M, N]^\tau - [M^\tau, N] &= [M^\tau, N^\tau] - [M^\tau, N] \\ &= \frac{1}{4} ([M^\tau + N^\tau] - [M^\tau - N^\tau]) - \frac{1}{4} ([M^\tau + N] - [M^\tau - N]) \\ &= \frac{1}{4} ([M^\tau + N^\tau] + [M^\tau - N]) - \frac{1}{4} ([M^\tau + N] + [M^\tau - N^\tau]), \end{aligned} \quad (32.12)$$

the difference of two increasing adapted processes. Also, this equals

$$\text{local martingale} - (M^\tau, N) + (M^\tau, N^\tau)$$

Claim: $(M^\tau, N) - (M^\tau, N^\tau) = (M^\tau, N - N^\tau)$ is a local martingale. Let σ_n be a localizing sequence for both M and N . Such a localizing sequence is of the form $\tau_n^M \wedge \tau_n^N$ where these are localizing sequences for the indicated local submartingale. Then obviously,

$$(-(M^\tau, N) + (M^\tau, N^\tau))^{\sigma_n} = -(M^{\sigma_n \wedge \tau}, N^{\sigma_n}) + (M^{\sigma_n \wedge \tau}, N^{\sigma_n \wedge \tau})$$

where N^{σ_n} and M^{σ_n} are martingales. To save notation, denote these by M and N respectively. Now use Lemma 32.1.1. Let σ be a stopping time with two values.

$$E((M^\tau(\sigma), N(\sigma) - N^\tau(\sigma))) = E(E((M^\tau(\sigma), N(\sigma) - N^\tau(\sigma)) | \mathcal{F}_\tau))$$

Now $M^\tau(\sigma)$ is $M(\sigma \wedge \tau)$ which is \mathcal{F}_τ measurable and so by the Doob optional sampling theorem,

$$\begin{aligned} &= E(M^\tau(\sigma), E(N(\sigma) - N^\tau(\sigma) | \mathcal{F}_\tau)) \\ &= E(M^\tau(\sigma), N(\sigma \wedge \tau) - N(\tau \wedge \sigma)) = 0 \end{aligned}$$

while

$$E((M^\tau(t), N(t) - N^\tau(t))) = E(E((M^\tau(t), N(t) - N^\tau(t)) | \mathcal{F}_\tau))$$

Since $M^\tau(t)$ is \mathcal{F}_τ measurable,

$$\begin{aligned} &= E((M^\tau(t), E(N(t) - N^\tau(t) | \mathcal{F}_\tau))) \\ &= E((M^\tau(t), E(N(t \wedge \tau) - N(t \wedge \tau)))) = 0 \end{aligned}$$

This shows the claim is true.

Now from 32.12 and Corollary 32.4.3, $[M, N]^\tau - [M^\tau, N] = 0$. Similarly $[M, N]^\tau - [M, N^\tau] = 0$. Now consider the next claim that $[M - M^\tau] = [M] - [M^\tau]$. From the definition, it follows

$$\begin{aligned} & [M - M^\tau] - ([M] + [M^\tau] - 2[M, M^\tau]) \\ &= \|M - M^\tau\|^2 - \left(\|M\|^2 + \|M^\tau\|^2 - 2(M, M^\tau) \right) + \text{local martingale} \\ &= \text{local martingale.} \end{aligned}$$

By the first part of the corollary which ensures $[M, M^\tau]$ is of bounded variation, the left side is the difference of two increasing adapted processes and so by Corollary 32.3.5 again, the left side equals 0. Thus from the above,

$$\begin{aligned} [M - M^\tau] &= [M] + [M^\tau] - 2[M, M^\tau] = [M] + [M^\tau] - 2[M^\tau, M^\tau] \\ &= [M] + [M^\tau] - 2[M^\tau] = [M] - [M^\tau] \leq [M] \end{aligned}$$

Finally consider the claim that $[M, N]$ is bilinear. From the definition, letting M_1, M_2, N be H valued local martingales,

$$\begin{aligned} (aM_1 + bM_2, N)_H &= [aM_1 + bM_2, N] + \text{local martingale} \\ a(M_1, N) + b(M_2, N)_H &= a[M_1, N] + b[M_2, N] + \text{local martingale} \end{aligned}$$

Hence

$$[aM_1 + bM_2, N] - (a[M_1, N] + b[M_2, N]) = \text{local martingale.}$$

The left side can be written as the difference of two increasing functions thanks to $[M, N]$ of bounded variation and so by Lemma 32.2.1 it equals 0. $[M, N]$ is obviously symmetric from the definition. ■

32.5 The Burkholder Davis Gundy Inequality

Define

$$M^*(\omega) \equiv \sup \{ \|M(t)(\omega)\| : t \in [0, T] \}.$$

The Burkholder Davis Gundy inequality is an amazing inequality which involves M^* and $[M](T)$.

Before presenting this, here is the good lambda inequality, Theorem 10.12.1 on Page 299 listed here for convenience.

Theorem 32.5.1 *Let $(\Omega, \mathcal{F}, \mu)$ be a finite measure space and let F be a continuous increasing function defined on $[0, \infty)$ such that $F(0) = 0$. Suppose also that for all $\alpha > 1$, there exists a constant C_α such that for all $x \in [0, \infty)$,*

$$F(\alpha x) \leq C_\alpha F(x).$$

Also suppose f, g are nonnegative measurable functions and there exists $\beta > 1, 0 < r \leq 1$, such that for all $\lambda > 0$ and $1 > \delta > 0$,

$$\mu([f > \beta\lambda] \cap [g \leq r\delta\lambda]) \leq \phi(\delta) \mu([f > \lambda]) \quad (32.13)$$

where $\lim_{\delta \rightarrow 0+} \phi(\delta) = 0$ and ϕ is increasing. Under these conditions, there exists a constant C depending only on β, ϕ, r such that

$$\int_{\Omega} F(f(\omega)) d\mu(\omega) \leq C \int_{\Omega} F(g(\omega)) d\mu(\omega).$$

The proof of the Burkholder Davis Gundy inequality also will depend on the hitting this before that theorem which is listed next for convenience.

Theorem 32.5.2 *Let $\{M(t)\}$ be a continuous real valued martingale adapted to the normal filtration \mathcal{F}_t and let*

$$M^* \equiv \sup \{|M(t)| : t \geq 0\}$$

and $M(0) = 0$. Letting

$$\tau_x \equiv \inf \{t > 0 : M(t) = x\}$$

Then if $a < 0 < b$ the following inequalities hold.

$$(b-a)P([\tau_b \leq \tau_a]) \geq -aP([M^* > 0]) \geq (b-a)P([\tau_b < \tau_a])$$

and

$$(b-a)P([\tau_a < \tau_b]) \leq bP([M^* > 0]) \leq (b-a)P([\tau_a \leq \tau_b]).$$

In words, $P([\tau_b \leq \tau_a])$ is the probability that $M(t)$ hits b no later than when it hits a . (Note that if $\tau_a = \infty = \tau_b$ then you would have $[\tau_a = \tau_b]$.)

Then the Burkholder Davis Gundy inequality is as follows. Generalizations will be presented later.

Theorem 32.5.3 *Let $\{M(t)\}$ be a continuous H valued martingale which is uniformly bounded, $M(0) = 0$, where H is a separable Hilbert space and $t \in [0, T]$. Then if F is a function of the sort described in the good lambda inequality above, there are constants, C and c independent of such martingales M such that*

$$c \int_{\Omega} F\left(\left([M](T)\right)^{1/2}\right) dP \leq \int_{\Omega} F(M^*) dP \leq C \int_{\Omega} F\left(\left([M](T)\right)^{1/2}\right) dP$$

where

$$M^*(\omega) \equiv \sup \{\|M(t)(\omega)\| : t \in [0, T]\}.$$

Proof: Using Corollary 32.4.3, let

$$\begin{aligned} N(t) &\equiv \|M(t) - M^{\tau}(t)\|^2 - [M - M^{\tau}](t) \\ &= \|M(t) - M^{\tau}(t)\|^2 - [M](t) + [M]^{\tau}(t) \end{aligned}$$

where

$$\tau \equiv \inf \{t \in [0, T] : \|M(t)\| > \lambda\}$$

Thus N is a martingale and $N(0) = 0$. In fact $N(t) = 0$ as long as $t \leq \tau$. As usual $\inf(\emptyset) \equiv \infty$. Note

$$[\tau < \infty] = \overset{\text{for some } t < T, \|M(t)\| > \lambda}{[M^* > \lambda]} \supseteq [N^* > 0]$$

This is because to say $\tau < \infty$ is to say there exists $t < T$ such that $\|M(t)\| > \lambda$ which is the same as saying $M^* > \lambda$. Thus the first two sets are equal. Either $\tau < \infty$ or $\tau = \infty$. If $\tau = \infty$, then from the formula for $N(t)$ above, $N(t) = 0$ for all $t \in [0, T]$ and so it can't happen that $N^* > 0$. Thus $[\tau = \infty] \subseteq [N^* = 0]$ so $[N^* > 0] \subseteq [\tau < \infty]$.

Let $\beta > 2$ and let $\delta \in (0, 1)$. Then $\beta - 1 > 1 > \delta > 0$. Consider the following which is set up to use the good lambda inequality.

$$S_r \equiv [M^* > \beta\lambda] \cap \left[([M](T))^{1/2} \leq r\delta\lambda \right]$$

where $0 < r < 1$. It is shown that S_r corresponds to hitting “this before that” and there is an estimate for this which involves $P([N^* > 0])$ which is bounded above by $P([M^* > \lambda])$ as discussed above. This will satisfy the hypotheses of the good lambda inequality.

Claim: For $\omega \in S_r$, $N(t)$ hits $\lambda^2(1 - \delta^2)$.

Proof of claim: For $\omega \in S_r$, there exists a $t < T$ such that $\|M(t)\| > \beta\lambda$ and so using Corollary 32.4.3 and triangle inequality,

$$\begin{aligned} N(t) &\geq \|M(t)\| - \|M^\tau(t)\|^2 - [M - M^\tau](t) \geq |\beta\lambda - \lambda|^2 - [M](t) \\ &\geq (\beta - 1)^2 \lambda^2 - \delta^2 \lambda^2 \end{aligned}$$

which shows that $N(t)$ hits $(\beta - 1)^2 \lambda^2 - \delta^2 \lambda^2$ for $\omega \in S_r$. By the intermediate value theorem, it also hits $\lambda^2(1 - \delta^2)$. This proves the claim.

Claim: $N(t)(\omega)$ never hits $-\delta^2 \lambda^2$ for $\omega \in S_r$.

Proof of claim: Suppose t is the first time $N(t)$ reaches $-\delta^2 \lambda^2$. Then $t > \tau$ because $N(t) = 0$ on $[0, \tau]$ and so

$$\begin{aligned} N(t) &= -\delta^2 \lambda^2 \geq \|M(t)\| - \lambda^2 - [M](t) + [M^\tau](t) \\ &\geq -r^2 \lambda^2 \delta^2, \end{aligned}$$

a contradiction since $r < 1$. This proves the claim.

Therefore, for all $\omega \in S_r$, $N(t)(\omega)$ reaches $\lambda^2(1 - \delta^2)$ before it reaches $-\delta^2 \lambda^2$. It follows

$$P(S_r) \leq P\left(N(t) \text{ reaches } \lambda^2(1 - \delta^2) \text{ before } -\delta^2 \lambda^2\right)$$

and because of Theorem 31.6.3 this is no larger than

$$P([N^* > 0]) \frac{\delta^2 \lambda^2}{\lambda^2(1 - \delta^2) - (-\delta^2 \lambda^2)} = P([N^* > 0]) \delta^2 \leq \delta^2 P([M^* > \lambda]).$$

Thus

$$P\left([M^* > \beta\lambda] \cap \left[([M](T))^{1/2} \leq r\delta\lambda\right]\right) \leq P([M^* > \lambda]) \delta^2$$

By the good lambda inequality,

$$\int_{\Omega} F(M^*) dP \leq C \int_{\Omega} F\left([M](T)^{1/2}\right) dP$$

which is one half the inequality.

Now consider the other half. This time define the stopping time τ by

$$\tau \equiv \inf\left\{t \in [0, T] : ([M](t))^{1/2} > \lambda\right\}$$

and let

$$S_r \equiv \left[([M](T))^{1/2} > \beta\lambda\right] \cap [2M^* \leq r\delta\lambda].$$

Then there exists $t < T$ such that $[M](t) > \beta^2 \lambda^2$. This time, let

$$N(t) \equiv [M](t) - [M^\tau](t) - \|M(t) - M^\tau(t)\|^2$$

This is still a martingale since by Corollary 32.4.3

$$[M](t) - [M^\tau](t) = [M - M^\tau](t)$$

Claim: $N(t)(\omega)$ hits $\lambda^2(1 - \delta^2)$ for some $t < T$ for $\omega \in S_r$.

Proof of claim: Fix such a $\omega \in S_r$. Let $t < T$ be such that $[M](t) > \beta^2 \lambda^2$. Then, since $\beta > 2$, $t > \tau$ and so for that ω ,

$$\begin{aligned} N(t) &> \beta^2 \lambda^2 - \lambda^2 - \|M(t) - M(\tau)\|^2 \\ &\geq (\beta - 1)^2 \lambda^2 - (\|M(t)\| + \|M(\tau)\|)^2 \\ &\geq (\beta - 1)^2 \lambda^2 - r^2 \delta^2 \lambda^2 \geq \lambda^2 - \delta^2 \lambda^2 \end{aligned}$$

By the intermediate value theorem, it hits $\lambda^2(1 - \delta^2)$. The last inequality follows because it is assumed that $2M^* \leq r\delta\lambda$. This proves the claim.

Claim: $N(t)(\omega)$ never hits $-\delta^2 \lambda^2$ for $\omega \in S_r$.

Proof of claim: By Corollary 32.4.3, if it did at t , then $t > \tau$ because $N(t) = 0$ for $t \leq \tau$, and so

$$\begin{aligned} 0 &\leq [M](t) - [M^\tau](t) = \|M(t) - M(\tau)\|^2 - \delta^2 \lambda^2 \\ &\leq (\|M(t)\| + \|M(\tau)\|)^2 - \delta^2 \lambda^2 \leq r^2 \delta^2 \lambda^2 - \delta^2 \lambda^2 < 0, \end{aligned}$$

a contradiction. The last inequality follows from $2M^* \leq r\delta\lambda$ on S_r . This proves the claim.

It follows that for each $r \in (0, 1)$,

$$P(S_r) \leq P\left(N(t) \text{ hits } \lambda^2(1 - \delta^2) \text{ before } -\delta^2 \lambda^2\right)$$

By Theorem 31.6.3 this is no larger than

$$\begin{aligned} P([N^* > 0]) \frac{\delta^2 \lambda^2}{\lambda^2(1 - \delta^2) + \delta^2 \lambda^2} &= P([N^* > 0]) \delta^2 \\ &\leq P([\tau < \infty]) \delta^2 = P\left(\left([M](T)\right)^{1/2} > \lambda\right) \delta^2 \end{aligned}$$

Now by the good lambda inequality, there is a constant k independent of M such that

$$\int_{\Omega} F\left(\left([M](T)\right)^{1/2}\right) dP \leq k \int_{\Omega} F(2M^*) dP \leq kC_2 \int_{\Omega} F(M^*) dP$$

by the assumptions about F . Therefore, combining this result with the first part,

$$\begin{aligned} (kC_2)^{-1} \int_{\Omega} F\left(\left([M](T)\right)^{1/2}\right) dP &\leq \int_{\Omega} F(M^*) dP \\ &\leq C \int_{\Omega} F\left(\left([M](T)\right)^{1/2}\right) dP \blacksquare \end{aligned}$$

Of course, everything holds for local martingales in place of martingales.

Theorem 32.5.4 Let $\{M(t)\}$ be a continuous H valued local martingale, $M(0) = 0$, where H is a separable Hilbert space and $t \in [0, T]$. Then if F is a function of the sort described in the good lambda inequality, that is,

$$F(0) = 0, F \text{ continuous, } F \text{ increasing,}$$

$$F(\alpha x) \leq c_\alpha F(x),$$

there are constants, C and c independent of such local martingales M such that

$$c \int_{\Omega} F([M](T)^{1/2}) dP \leq \int_{\Omega} F(M^*) dP \leq C \int_{\Omega} F([M](T)^{1/2}) dP$$

where

$$M^*(\omega) \equiv \sup \{\|M(t)(\omega)\| : t \in [0, T]\}.$$

Proof: Let $\{\tau_n\}$ be an increasing localizing sequence for M such that M^{τ_n} is uniformly bounded. Such a localizing sequence exists from Proposition 32.3.3. Then from Theorem 32.5.3 there exist constants c, C independent of τ_n such that

$$\begin{aligned} c \int_{\Omega} F([M^{\tau_n}](T)^{1/2}) dP &\leq \int_{\Omega} F((M^{\tau_n})^*) dP \\ &\leq C \int_{\Omega} F([M^{\tau_n}](T)^{1/2}) dP \end{aligned}$$

By Corollary 32.4.3, this implies

$$\begin{aligned} c \int_{\Omega} F([M^{\tau_n}](T)^{1/2}) dP &\leq \int_{\Omega} F((M^{\tau_n})^*) dP \\ &\leq C \int_{\Omega} F([M^{\tau_n}](T)^{1/2}) dP \end{aligned}$$

and now note that $[M^{\tau_n}](T)^{1/2}$ and $(M^{\tau_n})^*$ increase in n to $[M](T)^{1/2}$ and M^* respectively. Then the result follows from the monotone convergence theorem. ■

Here is a corollary [46].

Corollary 32.5.5 Let $\{M(t)\}$ be a continuous H valued local martingale and let $\varepsilon, \delta \in (0, \infty)$. Then there is a constant C , independent of ε, δ such that

$$P\left(\left[\sup_{t \in [0, T]} \|M(t)\| \geq \varepsilon\right]\right) \leq \frac{C}{\varepsilon} E\left([M]^{1/2}(T) \wedge \delta\right) + P\left([M]^{1/2}(T) > \delta\right)$$

Proof: Let the stopping time τ be defined by

$$\tau \equiv \inf \left\{ t > 0 : [M]^{1/2}(t) > \delta \right\}$$

Then

$$P([M^*] \geq \varepsilon) = P([M^*] \geq \varepsilon \cap [\tau = \infty]) + P([M^*] \geq \varepsilon \cap [\tau < \infty])$$

On the set where $[\tau = \infty]$, $M^\tau = M$ and so $P([M^*] \geq \varepsilon) \leq$

$$\leq \frac{1}{\varepsilon} \int_{\Omega} (M^\tau)^* dP + P\left([M^*] \geq \varepsilon \cap \left[[M]^{1/2}(T) > \delta\right]\right)$$

By Theorem 32.5.4 and Corollary 32.4.3,

$$\begin{aligned}
 &\leq \frac{C}{\varepsilon} \int_{\Omega} [M^c]^{1/2}(T) dP + P\left([M^* \geq \varepsilon] \cap \left[[M]^{1/2}(T) > \delta\right]\right) \\
 &= \frac{C}{\varepsilon} \int_{\Omega} ([M]^c)^{1/2}(T) dP + P\left([M^* \geq \varepsilon] \cap \left[[M]^{1/2}(T) > \delta\right]\right) \\
 &\leq \frac{C}{\varepsilon} \int_{\Omega} [M]^{1/2}(T) \wedge \delta dP + P\left([M^* \geq \varepsilon] \cap \left[[M]^{1/2}(T) > \delta\right]\right) \\
 &\leq \frac{C}{\varepsilon} \int_{\Omega} [M]^{1/2}(T) \wedge \delta dP + P\left(\left[[M]^{1/2}(T) > \delta\right]\right) \blacksquare
 \end{aligned}$$

The Burkholder Davis Gundy inequality along with the properties of the covariation implies the following amazing proposition.

Proposition 32.5.6 *The space $M_T^2(H)$ is a Hilbert space with respect to an equivalent norm. Here H is a separable Hilbert space.*

Proof: We already know from Proposition 31.7.2 that this space is a Banach space. It is only necessary to exhibit an equivalent norm which makes it a Hilbert space. However, you can let $F(\lambda) = \lambda^2$ in the Burkholder Davis Gundy theorem and obtain for $M \in M_T^2(H)$, the two norms

$$\left(\int_{\Omega} [M](T) dP\right)^{1/2} = \left(\int_{\Omega} [M, M](T) dP\right)^{1/2}$$

and

$$\left(\int_{\Omega} (M^*)^2 dP\right)^{1/2}$$

are equivalent. The first comes from an inner product since from Corollary 32.4.3, $[\cdot, \cdot]$ is bilinear and symmetric and nonnegative. If $[M, M](T) = [M](T) = 0$ in $L^1(\Omega)$, then from the Burkholder Davis Gundy inequality, $M^* = 0$ in $L^2(\Omega)$ and so $M = 0$. Hence

$$\int_{\Omega} [M, N](T) dP$$

is an inner product which yields the equivalent norm. \blacksquare

Later, the Wiener process will be discussed and the existence of such a process is proved. For now, the following example shows something about such processes.

Example 32.5.7 *An example of a real martingale is the Wiener process $W(t)$. It has the property that whenever $t_1 < t_2 < \dots < t_n$, the increments $\{W(t_i) - W(t_{i-1})\}$ are independent and whenever $s < t$, $W(t) - W(s)$ is normally distributed with mean 0 and variance $(t - s)$. For the Wiener process, we let*

$$\mathcal{F}_t \equiv \cap_{u>t} \overline{\sigma(W(s) - W(r) : r < s \leq u)}$$

and it is with respect to this normal filtration that W is a continuous martingale. What is the quadratic variation of such a process?

The quadratic variation of the Wiener process is just t . This is because if $A \in \mathcal{F}_s, s < t$,

$$\begin{aligned} E \left(\mathcal{X}_A \left(|W(t)|^2 - t \right) \right) &= \\ E \left(\mathcal{X}_A \left(|W(t) - W(s)|^2 + |W(s)|^2 + 2(W(s), W(t) - W(s)) - (t - s + s) \right) \right) \end{aligned}$$

Now

$$E \left(\mathcal{X}_A (2(W(s), W(t) - W(s))) \right) = P(A) E(2W(s)) E(W(t) - W(s)) = 0$$

by the independence of the increments. Thus the above reduces to

$$\begin{aligned} E \left(\mathcal{X}_A \left(|W(t) - W(s)|^2 + |W(s)|^2 - (t - s + s) \right) \right) \\ = E \left(\mathcal{X}_A \left(|W(t) - W(s)|^2 - (t - s) \right) \right) + E \left(\mathcal{X}_A \left(|W(s)|^2 - s \right) \right) \\ = P(A) E \left(|W(t) - W(s)|^2 - (t - s) \right) + E \left(\mathcal{X}_A \left(|W(s)|^2 - s \right) \right) \\ = E \left(\mathcal{X}_A \left(|W(s)|^2 - s \right) \right) \end{aligned}$$

and so $E \left(|W(t)|^2 - t | \mathcal{F}_s \right) = |W(s)|^2 - s$ showing that $t \rightarrow |W(t)|^2 - t$ is a martingale. Hence, by uniqueness, $[W](t) = t$.

32.6 Approximation With Step Functions

There is a really nice result about approximating a function $f \in L^p([0, T], E)$ with step functions. In this we deal with a specific representative of the equivalence class for $f \in L^p([0, T], E)$.

Lemma 32.6.1 *Let $f \in L^2([0, T]; E)$ for E a Banach space. For simplicity let f be Borel measurable. Then there exists a sequence of nested partitions, $\mathcal{P}_k \subseteq \mathcal{P}_{k+1}$,*

$$\mathcal{P}_k \equiv \{t_0^k, \dots, t_{m_k}^k\}$$

such that the step functions given by

$$\begin{aligned} f_k^r(t) &\equiv \sum_{j=1}^{m_k} f(t_j^k) \mathcal{X}_{[t_{j-1}^k, t_j^k)}(t) \\ f_k^l(t) &\equiv \sum_{j=1}^{m_k} f(t_{j-1}^k) \mathcal{X}_{[t_{j-1}^k, t_j^k)}(t) \end{aligned}$$

both converge to f in $L^2([0, T]; E)$ as $k \rightarrow \infty$ and

$$\lim_{k \rightarrow \infty} \max \left\{ |t_j^k - t_{j+1}^k| : j \in \{0, \dots, m_k\} \right\} = 0.$$

The mesh points $\left\{ t_j^k \right\}_{j=0}^{m_k}$ can be chosen to miss a given set of measure zero N if N does not contain either 0 or T .

Note that it would make no difference in terms of the conclusion of this lemma if you defined

$$f_k^l(t) \equiv \sum_{j=1}^{m_k} f\left(t_{j-1}^k\right) \mathcal{X}_{(t_{j-1}^k, t_j^k]}(t)$$

because the modified function equals the one given above off a countable subset of $[0, T]$, the union of the mesh points. One could change f_k^r similarly with no change in the conclusion.

Proof: Let f be 0 off $(0, T)$. Thus we will let it be defined on all of \mathbb{R} . Let $\gamma_n(t) \equiv k/2^n$, $\delta_n(t) \equiv (k+1)/2^n$, where $t \in (k/2^n, (k+1)/2^n]$, and $2^{-n} < \delta$. Thus $\gamma_n(t)$ is the closest $k2^{-n}$, $k \in \mathbb{Z}$ which is smaller than or equal to t while $\delta_n(t)$ is the closest $k2^{-n}$ larger than or equal to t . Let $g \in C_c((0, T), E)$ so the support of g is in $[\delta, T - \delta]$ for some $\delta > 0$ such that $\int_{[\delta, T-\delta]^c} \|f\|^p dt < \varepsilon$ and also

$$\int_0^T \|f - g\|^p dt = \int_{\mathbb{R}} \|f - g\|^p dt < \varepsilon.$$

Then

$$\begin{aligned} & \int_0^T \int_0^T \|f(\gamma_n(u-s) + s) - g(\gamma_n(u-s) + s)\|_E^p ds du \\ &= \int_0^T \int_0^T \|f(\gamma_n(u-s) + s) - g(\gamma_n(u-s) + s)\|_E^p duds \\ &= \int_0^T \int_{-s}^{T-s} \|f(\gamma_n(t) + s) - g(\gamma_n(t) + s)\|_E^p dt ds \\ &\leq \int_0^T \int_{-2T}^{2T} \|f(\gamma_n(t) + s) - g(\gamma_n(t) + s)\|_E^p dt ds \\ &= \int_{-2T}^{2T} \int_0^T \|f(\gamma_n(t) + s) - g(\gamma_n(t) + s)\|_E^p ds dt \\ &\leq \int_{-2T}^{2T} \int_{\mathbb{R}} \|f(\gamma_n(t) + s) - g(\gamma_n(t) + s)\|_E^p ds dt < 5T\varepsilon \end{aligned}$$

No effort is made to get the best possible estimate in the above. Then

$$\begin{aligned} & \int_0^T \int_0^T \|f(\gamma_n(u-s) + s) - f(u) - (g(\gamma_n(u-s) + s) - g(u))\|_E^p duds \\ &\leq 2^{p-1} \int_0^T \int_0^T \left(\|f(\gamma_n(u-s) + s) - g(\gamma_n(u-s) + s)\|_E^p + \|f(u) - g(u)\|_E^p \right) duds \\ &\leq 2^{p-1} \varepsilon 5T + 2^{p-1} \varepsilon T = 2^{p-1} 6\varepsilon T \end{aligned} \tag{32.14}$$

It follows that if n is chosen large enough, then

$$\int_0^T \int_0^T \|g(\gamma_n(u-s) + s) - g(u)\|_E^p duds < \varepsilon T$$

from uniform continuity of g . Therefore, from 32.14

$$\begin{aligned}
& \int_0^T \int_0^T \|f(\gamma_n(u-s) + s) - f(u)\|^p duds \\
& \leq 2^{p-1} \int_0^T \int_0^T \|g(\gamma_n(u-s) + s) - g(u)\|^p duds \\
& + 2^{p-1} \int_0^T \int_0^T \|f(\gamma_n(u-s) + s) - f(u) - (g(\gamma_n(u-s) + s) - g(u))\|_E^p duds \\
& \leq 2^{p-1} \varepsilon T + 2^{p-1} (2^{p-1} 6\varepsilon T) \tag{32.15}
\end{aligned}$$

Now $\gamma_n(u-s) + s \geq 0$ unless $u < \delta$. So consider

$$\begin{aligned}
& \int_0^T \int_0^T \|f((\gamma_n(u-s) + s)^+) - f(u)\|^p duds \leq \\
& \int_0^T \int_{2^{-n}}^T \|f(\gamma_n(u-s) + s) - f(u)\|^p duds + \int_0^T \int_0^{2^{-n}} \|f(u)\|^p duds \\
& \leq 2^{p-1} \varepsilon T + 2^{p-1} (2^{p-1} 6\varepsilon T) + \varepsilon \equiv \eta
\end{aligned}$$

Thus, since ε is arbitrary, so is η .

The function $u \rightarrow (\gamma_n(u-s) + s)^+$ has jumps $0 = t_0, t_1, t_2, \dots, t_{m_n-1}, t_{m_n} = T$ where these are listed in increasing order. The possible values of these t_i are $k2^{-n} + s$ for some $k \in \mathbb{Z}$. They are equally spaced being 2^{-n} apart except for the first two and the last two which are no more than 2^{-n} . One can slide this list of partition points around according to the choice of $s \in [0, T]$. Now suppose you have a set of measure zero N . Pick $s \in [0, T]$ such that none of the t_i are in N and

$$\int_0^T \|f((\gamma_n(u-s) + s)^+) - f(u)\|_E^p du < 2\eta$$

Now let $f_n(u) \equiv \sum_{k=1}^{m_n} f(t_{k-1}) \mathcal{X}_{[t_{k-1}, t_k)}(u)$. This is a step function of the desired sort. Then

$$\begin{aligned}
\int_0^T \|f_n(u) - f(u)\|_E^p du &= \sum_{k=1}^{m_n} \int_{t_{k-1}}^{t_k} \|f(t_{k-1}) - f(u)\|_E^p du \\
&= \sum_{k=1}^{m_n} \int_{t_{k-1}}^{t_k} \|f((\gamma_n(u-s) + s)^+) - f(u)\|_E^p du \\
&= \int_0^T \|f((\gamma_n(u-s) + s)^+) - f(u)\|_E^p du < 2\eta
\end{aligned}$$

Picking a sequence of these step functions f_j corresponding to $\eta = 2^{-j}$, one obtains the desired sequence in which values of f are assigned at the left end point of the interval. Making n still larger in the above argument and using the same argument with the right end points, one can also obtain a similar step function in which the values of f are given at the right end point which also converges to f in $L^2([0, T], E)$. ■

Chapter 33

Quadratic Variation and Stochastic Integration

Let \mathcal{F}_t be a normal filtration and let $\{M(t)\}$ be a continuous local martingale adapted to \mathcal{F}_t having values in U a separable real Hilbert space.

Definition 33.0.1 Let \mathcal{F}_t be a normal filtration and let

$$f(t) \equiv \sum_{k=0}^{m_n-1} f_k \mathcal{X}_{(t_k, t_{k+1}]}(t)$$

where $\{t_k\}_{k=0}^{m_n}$ is a partition of $[0, T]$ and each f_k is \mathcal{F}_{t_k} measurable, $f_k M^* \in L^2(\Omega)$ where

$$M^*(\omega) \equiv \sup_{t \in [0, T]} \|M(t)(\omega)\|$$

Such a function is called an elementary function. Also let $\{M(t)\}$ be a continuous local martingale adapted to \mathcal{F}_t which has values in a separable real Hilbert space U such that $M(0) = 0$. For such an elementary real valued function, define

$$\int_0^t f dM \equiv \sum_{k=0}^{m_n-1} f_k (M(t \wedge t_{k+1}) - M(t \wedge t_k)). \quad (33.1)$$

Since the $t \rightarrow \mathcal{F}_t$ is increasing, this definition is well defined. Also the set of elementary functions is a vector space.

Then with this definition, here is a wonderful lemma.

Lemma 33.0.2 For f an elementary function as above, $\{\int_0^t f dM\}$ is a continuous local martingale and

$$E \left(\left\| \int_0^t f dM \right\|_U^2 \right) = \int_{\Omega} \int_0^t f(s)^2 d[M](s) dP. \quad (33.2)$$

If N is another continuous local martingale adapted to \mathcal{F}_t and both f, g are elementary functions such that for each k ,

$$f_k M^*, g_k N^* \in L^2(\Omega),$$

then

$$E \left(\left(\int_0^t f dM, \int_0^t g dN \right)_U \right) = \int_{\Omega} \int_0^t f g d[M, N] \quad (33.3)$$

and both sides make sense.

Proof: Let $\{\tau_l\}$ be a localizing sequence for M such that M^{τ_l} is a bounded martingale. Then from the definition, for each ω

$$\int_0^t f dM = \lim_{l \rightarrow \infty} \int_0^t f dM^{\tau_l} = \lim_{l \rightarrow \infty} \left(\int_0^t f dM \right)^{\tau_l}$$

and it is clear that $\{\int_0^t f dM^{\tau_l}\}$ is a martingale because it is just the sum of some martingales. Thus $\{\tau_l\}$ is a localizing sequence for $\int_0^t f dM$. It is also clear $\int_0^t f dM$ is continuous

because it is a finite sum of continuous random variables. In the argument to get 33.2, I will write M rather than M^{τ_l} . Then it is understood that you can let $l \rightarrow \infty$ to obtain the desired formula for a local martingale. Thus, in what follows, M is a bounded martingale.

$$\begin{aligned} & E \left(\left\| \int_0^t f dM \right\|_U^2 \right) \\ &= E \left(\sum_{k=0}^{m_n-1} f_k (M(t \wedge t_{k+1}) - M(t \wedge t_k)), \sum_{j=0}^{m_n-1} f_j (M(t \wedge t_{j+1}) - M(t \wedge t_j)) \right) \end{aligned}$$

Let $M(t \wedge t_{k+1}) - M(t \wedge t_k) = \Delta M_k$. Thus ΔM_k is $\mathcal{F}_{t_{k+1}}$ measurable. Consider a mixed term in the above in which $j < k$

$$\begin{aligned} E(f_k \Delta M_k, f_j \Delta M_j) &= E(E(f_k \Delta M_k, f_j \Delta M_j) | \mathcal{F}_{t_k}) \\ &= E(f_k, f_j \Delta M_j E(\Delta M_k | \mathcal{F}_{t_k})) = 0 \end{aligned}$$

because $E(M(t \wedge t_{k+1}) - M(t \wedge t_k) | \mathcal{F}_{t_k}) = M(t \wedge t_{k+1} \wedge t_k) - M(t \wedge t_k)$. Thus

$$\begin{aligned} E \left(\left\| \int_0^t f dM \right\|_U^2 \right) &= \sum_{k=0}^{m_n-1} E(f_k \Delta M_k, f_k \Delta M_k) \\ &= \sum_{k=0}^{m_n-1} E(E((f_k \Delta M_k, f_k \Delta M_k) | \mathcal{F}_{t_k})) \end{aligned} \quad (33.4)$$

$$= \sum_{k=0}^{m_n-1} E(f_k^2 E((\Delta M_k, \Delta M_k) | \mathcal{F}_{t_k})) \quad (33.5)$$

now

$$\begin{aligned} E((M(t \wedge t_{k+1}), M(t \wedge t_k)) | \mathcal{F}_{t_k}) &= (E(M(t \wedge t_{k+1}) | \mathcal{F}_{t_k}), M(t \wedge t_k)) \\ &= \|M(t \wedge t_k)\|^2 \end{aligned}$$

and so

$$\begin{aligned} E((\Delta M_k, \Delta M_k) | \mathcal{F}_{t_k}) &= E \left(\begin{pmatrix} M(t \wedge t_{k+1}) - M(t \wedge t_k) \\ M(t \wedge t_{k+1}) - M(t \wedge t_k) \end{pmatrix} | \mathcal{F}_{t_k} \right) = \\ &= E \left(\|M(t \wedge t_{k+1})\|^2 | \mathcal{F}_{t_k} \right) + E \left(\|M(t \wedge t_k)\|^2 | \mathcal{F}_{t_k} \right) \\ &\quad - 2E((M(t \wedge t_{k+1}), M(t \wedge t_k)) | \mathcal{F}_{t_k}) \\ &= E \left(\|M(t \wedge t_{k+1})\|^2 | \mathcal{F}_{t_k} \right) - \|M(t \wedge t_k)\|^2 \end{aligned}$$

Therefore, the right side of 33.5 is

$$\sum_{k=0}^{m_n-1} E \left(f_k^2 \|M(t \wedge t_{k+1})\|^2 \right) - E \left(f_k^2 \|M(t \wedge t_k)\|^2 \right).$$

Now recall that $\|M(t)\|^2 = [M](t) + N(t)$ where $N(t)$ is a martingale. It then reduces to

$$\sum_{k=0}^{m_n-1} E \left(f_k^2 ([M](t \wedge t_{k+1}) + N(t \wedge t_{k+1})) \right) - E \left(f_k^2 ([M](t \wedge t_k) + N(t \wedge t_k)) \right)$$

$$= \int_{\Omega} \int_0^t f(s)^2 d[M] dP$$

since the martingales integrate to 0. This proves the first formula. If $f(s) \in \mathcal{L}(U, H)$, you could probably modify this argument. In this case, you couldn't factor out of the inner product because it would no longer be a scalar. However, if $f(s) \in \mathcal{L}_2(U, H)$ you maybe could do something. This is really a version of the Ito isometry discussed later.

Next is a similar argument for the case of two different elementary functions. There is no loss of generality in assuming the mesh points are the same for the two elementary functions because if not, one can simply add in points to make this happen. It suffices to consider 33.3 because the other formula is a special case. To begin with, let $\{\tau_l\}$ be a localizing sequence which makes both M^{τ_l} and N^{τ_l} into bounded martingales. Consider the stopped process.

$$\begin{aligned} & E \left(\left(\int_0^t f dM^{\tau_l}, \int_0^t g dN^{\tau_l} \right)_U \right) \\ &= E \left(\left(\sum_{k=0}^{m_n-1} f_k (M^{\tau_l}(t \wedge t_{k+1}) - M^{\tau_l}(t \wedge t_k)) , \right. \right. \\ & \quad \left. \left. \sum_{k=0}^{m_n-1} g_k (N^{\tau_l}(t \wedge t_{k+1}) - N^{\tau_l}(t \wedge t_k)) \right) \right) \end{aligned}$$

To save on notation, write $M^{\tau_l}(t \wedge t_{k+1}) - M^{\tau_l}(t \wedge t_k) \equiv \Delta M_k(t)$, similar for ΔN_k . Thus

$$\Delta M_k = M^{\tau_l \wedge t_{k+1}} - M^{\tau_l \wedge t_k},$$

similar for ΔN_k . Then the above equals

$$E \left(\sum_{k=0}^{m_n-1} \left(f_k \Delta M_k, \sum_{k=0}^{m_n-1} g_k \Delta N_k \right) \right) = E \left(\sum_{k,j} f_k g_j (\Delta M_k, \Delta N_j) \right)$$

Now consider one of the mixed terms with $j < k$.

$$\begin{aligned} E((f_k \Delta M_k, g_j \Delta N_j)) &= E(E((f_k \Delta M_k, g_j \Delta N_j) | \mathcal{F}_{t_k})) \\ &= E(g_j \Delta N_j, f_k E(\Delta M_k | \mathcal{F}_{t_k})) = 0 \end{aligned}$$

since $E(\Delta M_k | \mathcal{F}_{t_k}) = E((M^{\tau_l}(t \wedge t_{k+1}) - M^{\tau_l}(t \wedge t_k)) | \mathcal{F}_{t_k}) = 0$ by the Doob optional sampling theorem. Thus

$$E \left(\left(\int_0^t f dM^{\tau_l}, \int_0^t g dN^{\tau_l} \right)_U \right) = \quad (33.6)$$

$$= \sum_{k=0}^{m_n-1} E(f_k g_k (\Delta M_k, \Delta N_k)) = \sum_{k=0}^{m_n-1} E(f_k g_k ([\Delta M_k, \Delta N_k] + \mathcal{N}_k)) \quad (33.7)$$

where \mathcal{N}_k is a martingale such that $\mathcal{N}_k(t) = 0$ for all $t \leq t_k$. This is because the martingale $(N^{\tau_l})^{t_{k+1}} - (N^{\tau_l})^{t_k} = \Delta N_k$ equals 0 for such t ; and so $E(\mathcal{N}_k(t)) = 0$. Thus $f_k g_k \mathcal{N}_k$ is a martingale which equals zero when $t = 0$. Therefore, its expectation also equals 0. Consequently the above reduces to

$$\sum_{k=0}^{m_n-1} E(f_k g_k [\Delta M_k, \Delta N_k]).$$

At this point, recall the definition of the covariation. The above equals

$$\frac{1}{4} \sum_{k=0}^{m_n-1} E(f_k g_k ([\Delta M_k + \Delta N_k] - [\Delta M_k - \Delta N_k]))$$

Rewriting this yields

$$\begin{aligned} &= \frac{1}{4} \sum_{k=0}^{m_n-1} E(f_k g_k ([(M^{\tau_l})^{t_{k+1}} + (N^{\tau_l})^{t_{k+1}} - ((M^{\tau_l})^{t_k} + (N^{\tau_l})^{t_k})] \\ &\quad - [(M^{\tau_l})^{t_{k+1}} - (N^{\tau_l})^{t_{k+1}} - ((M^{\tau_l})^{t_k} - (N^{\tau_l})^{t_k})])) \end{aligned}$$

To save on notation, denote

$$\begin{aligned} (M^{\tau_l})^{t_{k+1}} + (N^{\tau_l})^{t_{k+1}} - ((M^{\tau_l})^{t_k} + (N^{\tau_l})^{t_k}) &\equiv \Delta_k (M^{\tau_l} + N^{\tau_l}) \\ (M^{\tau_l})^{t_{k+1}} - (N^{\tau_l})^{t_{k+1}} - ((M^{\tau_l})^{t_k} - (N^{\tau_l})^{t_k}) &\equiv \Delta_k (M^{\tau_l} - N^{\tau_l}) \end{aligned}$$

Thus the above equals

$$\frac{1}{4} \sum_{k=0}^{m_n-1} E(f_k g_k ([\Delta_k (M^{\tau_l} + N^{\tau_l})] - [\Delta_k (M^{\tau_l} - N^{\tau_l})]))$$

Now from Corollary 32.4.3,

$$= \frac{1}{4} \sum_{k=0}^{m_n-1} E(f_k g_k ([\Delta_k (M + N)]^{\tau_l} - [\Delta_k (M - N)]^{\tau_l}))$$

Letting $l \rightarrow \infty$, this reduces to

$$\begin{aligned} &= \frac{1}{4} \sum_{k=0}^{m_n-1} E(f_k g_k ([\Delta_k (M + N)] - [\Delta_k (M - N)])) \\ &= \frac{1}{4} \left(\int_{\Omega} \int_0^t f g (d[M + N] - d[M - N]) \right) \\ &= \int_{\Omega} \int_0^t f g d[M, N] \end{aligned}$$

Now consider the left side of 33.7.

$$\begin{aligned} &E \left(\left(\int_0^t f dM^{\tau_l}, \int_0^t g dN^{\tau_l} \right)_U \right) \\ &\equiv \int_{\Omega} \sum_{k,j} f_k g_j ((M^{\tau_l}(t \wedge t_{k+1}) - M^{\tau_l}(t \wedge t_k)), \\ &\quad (N^{\tau_l}(t \wedge t_{j+1}) - N^{\tau_l}(t \wedge t_j))) dP \end{aligned}$$

Then for each ω , the integrand converges as $l \rightarrow \infty$ to

$$\sum_{k,j} f_k g_j ((M(t \wedge t_{k+1}) - M(t \wedge t_k)), (N(t \wedge t_{j+1}) - N(t \wedge t_j)))$$

But also you can do a sloppy estimate which will allow the use of the dominated convergence theorem.

$$\left\| \sum_{k,j} f_k g_j (M^{\tau_l}(t \wedge t_{k+1}) - M^{\tau_l}(t \wedge t_k)), (N^{\tau_l}(t \wedge t_{j+1}) - N^{\tau_l}(t \wedge t_j)) \right\|$$

$$\leq \sum_{k,j} |f_k| |g_j| 4M^* N^* \in L^1(\Omega)$$

by assumption. Thus the left side of 33.7 converges as $l \rightarrow \infty$ to

$$\int_{\Omega} \sum_{k,j} f_k g_j ((M(t \wedge t_{k+1}) - M(t \wedge t_k)), (N(t \wedge t_{j+1}) - N(t \wedge t_j))) dP$$

$$= \int_{\Omega} \left(\int_0^t f dM, \int_0^t g dN \right)_U dP \blacksquare$$

Note for each ω , the inside integral in 33.2 is just a Stieltjes integral taken with respect to the increasing integrating function $[M]$.

Of course, with this estimate it is obvious how to extend the integral to a larger class of functions.

Definition 33.0.3 Let $\nu(\omega)$ denote the Radon measure representing the functional

$$\Lambda(\omega)(g) \equiv \int_0^T g d[M](t)(\omega)$$

($t \rightarrow [M](t)(\omega)$ is a continuous increasing function and $\nu(\omega)$ is the measure representing the Stieltjes integral, one for each ω .) Then let \mathcal{G}_M denote functions $f(s, \omega)$ which are the limit of such elementary functions in the space $L^2(\Omega; L^2([0, T], \nu(\cdot)))$, the norm of such functions being

$$\|f\|_{\mathcal{G}}^2 \equiv \int_{\Omega} \int_0^T f(s)^2 d[M](s) dP$$

For $f \in \mathcal{G}$ just defined,

$$\int_0^t f dM \equiv \lim_{n \rightarrow \infty} \int_0^t f_n dM$$

where $\{f_n\}$ is a sequence of elementary functions converging to f in

$$L^2(\Omega; L^2([0, T], \nu(\cdot))).$$

Now here is an interesting lemma.

Lemma 33.0.4 Let M, N be continuous local martingales, $M(0) = N(0) = 0$ having values in a separable Hilbert space, U . Then

$$[M + N]^{1/2} \leq ([M]^{1/2} + [N]^{1/2}) \quad (33.8)$$

$$[M + N] \leq 2([M] + [N]) \quad (33.9)$$

Also, let ν_{M+N} denote the measure obtained from the increasing function $[M + N]$ and ν_N, ν_M be defined similarly,

$$\nu_{M+N} \leq 2(\nu_M + \nu_N) \quad (33.10)$$

on all Borel sets.

Proof: Since $(M, N) \rightarrow [M, N]$ is bilinear and satisfies

$$\begin{aligned} [M, N] &= [N, M] \\ [aM + bM_1, N] &= a[M, N] + b[M_1, N] \\ [M, M] &\geq 0 \end{aligned}$$

which follows from Corollary 32.4.3, the usual Cauchy Schwarz inequality holds and so

$$|[M, N]| \leq [M]^{1/2} [N]^{1/2}$$

Thus

$$\begin{aligned} [M + N] &\equiv [M + N, M + N] = [M, M] + [N, N] + 2[M, N] \\ &\leq [M] + [N] + 2[M]^{1/2} [N]^{1/2} = \left([M]^{1/2} + [N]^{1/2}\right)^2 \end{aligned}$$

This proves 33.8. Now square both sides. Then the right side is no larger than

$$2([M] + [N])$$

and this shows 33.9.

Now consider the claim about the measures. It was just shown that

$$[(M + N) - (M + N)^s] \leq 2([M - M^s] + [N - N^s])$$

and from Corollary 32.4.3 this implies that for $t > s$

$$\begin{aligned} &[M + N](t) - [M + N](s \wedge t) \\ &= [M + N](t) - [M + N]^s(t) \\ &= [M + N - (M^s + N^s)](t) \\ &= [M - M^s + (N - N^s)](t) \\ &\leq 2[M - M^s](t) + 2[N - N^s](t) \\ &\leq 2([M](t) - [M](s)) + 2([N](t) - [N](s)) \end{aligned}$$

Thus

$$\nu_{M+N}([s, t]) \leq 2(\nu_M([s, t]) + \nu_N([s, t]))$$

By regularity of the measures, this continues to hold with any Borel set F in place of $[s, t]$.

■

Theorem 33.0.5 *The integral is well defined and has a continuous version which is a local martingale. Furthermore it satisfies the Ito isometry,*

$$E \left(\left\| \int_0^t f dM \right\|_U^2 \right) = \int_{\Omega} \int_0^t f(s)^2 d[M](s) dP \quad (33.11)$$

Let the norm on $\mathcal{G}_N \cap \mathcal{G}_M$ be the maximum of the norms on \mathcal{G}_N and \mathcal{G}_M and denote by \mathcal{E}_N and \mathcal{E}_M the elementary functions corresponding to the martingales N and M respectively. Define \mathcal{G}_{NM} as the closure in $\mathcal{G}_N \cap \mathcal{G}_M$ of $\mathcal{E}_N \cap \mathcal{E}_M$. Then for $f, g \in \mathcal{G}_{NM}$,

$$E \left(\left(\int_0^t f dM, \int_0^t g dN \right) \right) = \int_{\Omega} \int_0^t f g d[M, N] \quad (33.12)$$

Proof: It is clear the definition is well defined because if $\{f_n\}$ and $\{g_n\}$ are two sequences of elementary functions converging to f in $L^2(\Omega; L^2([0, T], \nu))$ and if $\int_0^t f dM$ is the integral which comes from $\{g_n\}$,

$$\begin{aligned} & \int_{\Omega} \left\| \int_0^t f dM - \int_0^t f_n dM \right\|^2 dP \\ &= \lim_{n \rightarrow \infty} \int_{\Omega} \left\| \int_0^t g_n dM - \int_0^t f_n dM \right\|^2 dP \\ &\leq \lim_{n \rightarrow \infty} \int_{\Omega} \int_0^T \|g_n - f_n\|^2 d\nu dP = 0. \end{aligned}$$

Consider the claim the integral has a continuous version. Recall Theorem 31.4.3, part of which is listed here for convenience.

Theorem 33.0.6 *Let $\{X(t)\}$ be a right continuous nonnegative submartingale adapted to the normal filtration \mathcal{F}_t for $t \in [0, T]$. Let $p \geq 1$. Define*

$$X^*(t) \equiv \sup \{X(s) : 0 < s < t\}, \quad X^*(0) \equiv 0.$$

Then for $\lambda > 0$

$$P([X^*(T) > \lambda]) \leq \frac{1}{\lambda^p} \int_{\Omega} X(T)^p dP \quad (33.13)$$

Let $\{f_n\}$ be a sequence of elementary functions converging to f in

$$L^2(\Omega; L^2([0, T], \nu(\cdot))).$$

Then letting

$$\begin{aligned} X_{n,m}^{\tau_l}(t) &= \left\| \int_0^t (f_n - f_m) dM^{\tau_l} \right\|_U, \\ X_{n,m}(t) &= \left\| \int_0^t (f_n - f_m) dM \right\|_U = \left\| \int_0^t f_n dM - \int_0^t f_m dM \right\|_U \end{aligned}$$

It follows $X_{n,m}^{\tau_l}$ is a continuous nonnegative submartingale and from Theorem 31.4.3 just listed,

$$\begin{aligned} P([X_{n,m}^{\tau_l*}(T) > \lambda]) &\leq \frac{1}{\lambda^2} \int_{\Omega} X_{n,m}^{\tau_l}(T)^2 dP \\ &\leq \frac{1}{\lambda^2} \int_{\Omega} \int_0^T |f_n - f_m|^2 d[M^{\tau_l}] dP \\ &\leq \frac{1}{\lambda^2} \int_{\Omega} \int_0^T |f_n - f_m|^2 d[M] dP \end{aligned}$$

Letting $l \rightarrow \infty$,

$$P([X_{n,m}^*(T) > \lambda]) \leq \frac{1}{\lambda^2} \int_{\Omega} \int_0^T |f_n - f_m|^2 d[M] dP$$

Therefore, there exists a subsequence, still denoted by $\{f_n\}$ such that

$$P([X_{n,n+1}^*(T) > 2^{-n}]) < 2^{-n}$$

Then by the Borel Cantelli lemma, the ω in infinitely many of the sets

$$[X_{n,n+1}^*(T) > 2^{-n}]$$

has measure 0. Denoting this exceptional set as N , it follows that for $\omega \notin N$, there exists $n(\omega)$ such that for $n > n(\omega)$,

$$\sup_{t \in [0, T]} \left\| \int_0^t f_n dM - \int_0^t f_{n+1} dM \right\| \leq 2^{-n}$$

and this implies uniform convergence of $\{\int_0^t f_n dM\}$. Letting

$$G(t) = \lim_{n \rightarrow \infty} \int_0^t f_n dM,$$

for $\omega \notin N$ and $G(t) = 0$ for $\omega \in N$, it follows that for each t , the continuous adapted process $G(t)$ equals $\int_0^t f dM$ a.e. Thus $\{\int_0^t f dM\}$ has a continuous version.

It suffices to verify 33.12. Let $\{f_n\}$ and $\{g_n\}$ be sequences of elementary functions converging to f and g in $\mathcal{G}_M \cap \mathcal{G}_N$. By Lemma 33.0.2,

$$E \left(\left(\int_0^t f_n dM, \int_0^t g_n dN \right)_U \right) = \int_{\Omega} \int_0^t f_n g_n d[M, N]$$

Then by the Holder inequality and the above definition,

$$\lim_{n \rightarrow \infty} E \left(\left(\int_0^t f_n dM, \int_0^t g_n dN \right)_U \right) = E \left(\left(\int_0^t f dM, \int_0^t g dN \right)_U \right)$$

Consider the right side which equals

$$\frac{1}{4} \int_{\Omega} \int_0^t f_n g_n d[M + N] dP - \frac{1}{4} \int_{\Omega} \int_0^t f_n g_n d[M - N] dP$$

Now from Lemma 33.0.4,

$$\begin{aligned} & \left| \int_{\Omega} \int_0^t f_n g_n d[M + N] dP - \int_{\Omega} \int_0^t f g d[M + N] dP \right| \\ &= \left| \int_{\Omega} \int_0^t f_n g_n d\mathbf{v}_{M+N} dP - \int_{\Omega} \int_0^t f g d\mathbf{v}_{M+N} dP \right| \\ &\leq 2 \left(\int_{\Omega} \int_0^t |f_n g_n - f g| d\mathbf{v}_M dP + \int_{\Omega} \int_0^t |f_n g_n - f g| d\mathbf{v}_N dP \right) \end{aligned}$$

and by the choice of the f_n and g_n , these both converge to 0. Similar considerations apply to

$$\left| \int_{\Omega} \int_0^t f_n g_n d[M - N] dP - \int_{\Omega} \int_0^t f g d[M - N] dP \right|$$

and show

$$\lim_{n \rightarrow \infty} \int_{\Omega} \int_0^t f_n g_n d[M, N] = \int_{\Omega} \int_0^t f g d[M, N] \blacksquare$$

33.1 The Stieltjes Integral

When we do Stieltjes integration, the first thing to consider is continuous functions. I am going to do this here. It seems to me that this is an important case to consider if for no other reason than this is what we do with Stieltjes integration. If M is continuous and f is of bounded variation, and $\{t_k^n\}_{k=0}^{m_n-1}$ is a partition P_n of $[0, T]$ for which $\|P_n\| \equiv \max \{|t_{k+1}^n - t_k^n|, k = 0, 1, \dots, m_n - 1\}$, then

$$\lim_{n \rightarrow \infty} \sum_{k=0}^{m_n-1} f(t_k^n) (M(t \wedge t_{k+1}^n) - M(t \wedge t_k^n)) = \int_0^t f dM \quad (33.14)$$

exists in U . This is because of the integration by parts theorem for Stieltjes integrals. Indeed, $\int_0^t M df$ exists by standard arguments. This is a little different here because M has values in a Hilbert space, but the proof is essentially the same as for scalar valued functions. See the proof in my single variable advanced calculus book, or [2] or [28]. Thus for a.e. ω , one obtains the limit in 33.14 directly from the theory of Stieltjes integration, this limit taking place in U . Suppose now that

$$M^* \in L^2(\Omega)$$

In the above context, 33.11 is

$$E \left(\left\| \int_0^t f dM \right\|_U^2 \right) = \int_{\Omega} \int_0^t f(s)^2 d[M](s) dP$$

when f is an elementary function. Let $f(t, \omega)$ be bounded and continuous in t with $f(t, \cdot)$ \mathcal{F}_t measurable. Let $P_n = \{t_j^n\}_{j=0}^{m_n-1}$ and let f_k be the elementary function

$$f_k(t, \omega) \equiv \sum_{j=0}^{m_n-1} f_k(t_j^n) \mathcal{X}_{(t_j^n, t_{j+1}^n]}(t)$$

Thus these elementary functions converge uniformly to f for $t \in [0, T]$ for fixed ω as $\|P_n\| \rightarrow 0$ and they are all bounded uniformly. Therefore, from 33.11 and the maximal estimates of Theorem 30.5.3,

$$\begin{aligned} & \frac{1}{2} E \left(\sup_{t \leq T} \left\| \int_0^t f_k dM - \int_0^t f_m dM \right\|_U^2 \right) \leq \\ & E \left(\left\| \int_0^T f_k dM - \int_0^T f_m dM \right\|_U^2 \right) = \int_{\Omega} \int_0^T |f_k(s) - f_m(s)|^2 d[M](s) dP \end{aligned} \quad (33.15)$$

This is because you can assume, by taking the union of the two partitions involved, that you are dealing with a single partition for both f_k and f_m . Now for a.e. ω ,

$$\lim_{k, m \rightarrow \infty} \int_0^T |f_k(s) - f_m(s)|^2 d[M](s) = 0.$$

This is because of the Burkholder Davis Gundy inequality which implies $d[M]$ is a finite measure for a.e. ω . Indeed, from this inequality,

$$c \int_{\Omega} \left(([M](T))^{1/2} \right)^2 dP = c \int_{\Omega} [M](T) dP$$

$$= c \int_{\Omega} \int_0^T d[M] dP \leq \int_{\Omega} (M^*)^2 dP < \infty \quad (33.16)$$

Hence for a.e. ω , $\lim_{k,m \rightarrow \infty} \int_0^T |f_k(s) - f_m(s)|^2 d[M](s) = 0$. Also from the fact these elementary functions are all bounded, that inside integral on the right in 33.15 is no more than $K[M(T)]$ for a constant K which comes from the upper bound of all these elementary functions. This is a function in $L^1(\Omega)$ by 33.16. Now if A is a measurable set,

$$\int_A \int_0^T |f_{n+k}(s) - f_n(s)|^2 d[M](s) dP \leq K \int_A [M](T) dP$$

and $[M](T)$ is a function in L^1 so this collection of functions of ω ,

$$\int_0^T |f_k(s) - f_m(s)|^2 d[M](s)$$

is uniformly integrable. By the Vitali convergence theorem,

$$\lim_{k,m \rightarrow \infty} \int_{\Omega} \int_0^T |f_k(s) - f_m(s)|^2 d[M](s) dP = 0.$$

By 33.15, $\{\int_0^t f_k dM\}_{k=1}^{\infty}$ is a Cauchy sequence in $\mathcal{M}_T^2(U)$ and so there is a unique martingale $I(t) \equiv \int_0^t f dM$ such that $\int_0^t f_k dM \rightarrow I(t)$ in $\mathcal{M}_T^2(U)$. Also by Proposition 31.7.2, there is a subsequence which converges uniformly on $[0, T]$ for a.e. ω . In case f is of bounded variation in addition to being continuous because in this case, the Stieltjes sums $\int_0^t f_{n_k} dM$ converge to $\int_0^t f dM$. This proves the following interesting relationship.

Proposition 33.1.1 *In the above context where $M^* \in L^2(\Omega)$, M a martingale with values in U a separable Hilbert space, suppose $t \rightarrow f(t, \omega)$ is of bounded variation and is continuous and $\omega \rightarrow f(t, \omega)$ is adapted to the filtration \mathcal{F}_t . Also suppose $(t, \omega) \rightarrow f(t, \omega)$ is bounded. Then the ordinary Stieltjes integral $\int_0^t f(s) dM(\omega)$ is a martingale.*

In the above argument, it was not necessary that $t \rightarrow f(t, \omega)$ have bounded variation so the $I(t)$ is in a sense more general than the Stieltjes integral but it extends the idea of the Stieltjes integral.

What if M is only a local martingale? Then you could let σ_m be the first hitting time of m by $\|M(t)\|$ and you could repeat everything and get

$$\int_0^t f(s) dM^{\sigma_m} = \int_0^{t \wedge \sigma_m} f(s) dM$$

is a martingale. The approximating sums in this case would be

$$\begin{aligned} & \sum_{k=0}^{m_n-1} f(t_k^n) (M^{\sigma_m}(t \wedge t_{k+1}^n) - M^{\sigma_m}(t \wedge t_k^n)) \rightarrow \int_0^t f(s) dM^{\sigma_m} \\ & \sum_{k=0}^{m_n-1} f(t_k^n) (M^{\sigma_m}(t \wedge t_{k+1}^n) - M^{\sigma_m}(t \wedge t_k^n)) \\ &= \left(\sum_{k=0}^{m_n-1} f(t_k^n) (M(t \wedge t_{k+1}^n) - M(t \wedge t_k^n)) \right)^{\sigma_m} \rightarrow \int_0^{t \wedge \sigma_m} f(s) dM \end{aligned}$$

as $\|P_n\| \rightarrow \infty$. Thus, more generally, this Stieltjes integral $\int_0^t f dM$ is a local martingale. This is all written with scalar valued $f(s, \omega)$ in mind, but if $f(s, \omega)$ were something in $\mathcal{L}(U, U)$ would it be any different? In case $[M]$ depends only on t , all the above considerations become easier. Indeed, if $[M]$ is just an increasing function of t , 33.15 would imply that one could get $\int_0^t f_k(s) dM$ is a Cauchy sequence in $\mathcal{M}_T^2(U)$ whenever $\lim_{k,m \rightarrow \infty} \int_0^T \int_0^T |f_k(s) - f_m(s)|^2 dF(s) dP = 0$ where $F(s) = [M](s)$ and F an increasing function. In fact, the main interest will be when $[M](t) = t$ as in Example 32.5.7 so $\{\int_0^t f_k(s) dM\}$ being a Cauchy sequence in $\mathcal{M}_T^2(U)$ comes from assuming simply that $\{f_k\}$ is Cauchy in $L^2(\Omega; L^2([0, T]))$.

33.2 The Stochastic Integral When $f(s) \in \mathcal{L}_2(U, H)$

Let H, U be separable Hilbert spaces and suppose for each $s, f(s) \in \mathcal{L}_2(U, H)$, the space of Hilbert Schmidt operators described in the section on compact operators Section 22.5.2. Recall that $f(s) \in \mathcal{L}_2(U, H)$ implies $f(s)^* f(s)$ is a self adjoint compact operator thanks to Theorem 22.5.18. Thus as pointed out there, we can pick **any** orthonormal basis $\{e_k\}$ for U and

$$\|f^*(s) f(s)\|_{\mathcal{L}_2}^2 = \sum_{k=1}^{\infty} \|f^*(s) f(s) e_k\|_H^2$$

the same value being obtained for any of these orthonormal sets. Now also $f^*(s) f(s)$ is compact and self adjoint so by the Hilbert Schmidt theorem, Theorem 22.5.3, there is a decreasing list of positive numbers $\{\lambda_k\}$ and a corresponding orthonormal set of eigenvectors $\{e_k\}$ such that $f^*(s) f(s) e_k = \lambda_k e_k$. Then from this equation, $\|f(s) e_k\|^2 = \lambda_k$ and so $\sum_k |\lambda_k| < \infty$.

Also, for $L \in \mathcal{L}_2(U, H)$ since $L^* L$ is nonnegative, there is a self adjoint square root $\sqrt{L^* L}$ and $\sqrt{L^* L} e_k = \sqrt{\lambda_k} e_k$. Note also, $\|\sqrt{L^* L} e_k\|^2 = \|L e_k\|^2$. As pointed out earlier, This implies that for any pair of orthonormal basis $\{e_k\}, \{\hat{e}_k\}$

$$\sum_k \|L e_k\|^2 = \sum_k \|\sqrt{L^* L} e_k\|^2 = \sum_k \|\sqrt{L^* L} \hat{e}_k\|^2 = \sum_k \|L \hat{e}_k\|^2$$

so the norm of something in $\mathcal{L}_2(U, H)$ can be defined using any orthonormal basis.

Definition 33.2.1 *An elementary function is one of the form*

$$f(t) \equiv \sum_{r=0}^{m_n-1} f_r \mathcal{X}_{(t_r, t_{r+1}]}(t)$$

where $f_r \in \mathcal{L}_2(U, H)$ and f_r is \mathcal{F}_{t_r} measurable. In this section, assume also that f_r is bounded. Here $t_0 < t_1 < \dots < t_{m_n}$ is a partition of $[0, T]$ as above. We can define the stochastic integral of an elementary function with respect to a continuous martingale $M(t)$ as

$$\int_0^t f dM \equiv \sum_{r=0}^{m_n-1} f_r (M(t \wedge t_{r+1}) - M(t \wedge t_r))$$

Note that $\int_0^t f dM \in H$ for a given ω . Also note that, since the \mathcal{F}_t are increasing, given two elementary functions, we can write them both with respect to the same partition and consequently this set of elementary functions is a linear space and the integral on these functions is linear.

Note that for $t_r < s \leq t$,

$$\begin{aligned} E(f_r(M(t \wedge t_{r+1}) - M(t \wedge t_r)) | \mathcal{F}_s) &= f_r E(M(t \wedge t_{r+1}) - M(t \wedge t_r) | \mathcal{F}_s) \\ &= f_r(M(s \wedge t_{r+1}) - M(s \wedge t_r)) \end{aligned}$$

However, for $t \leq t_r$ the term in the sum equals 0. Thus each term in that sum is a martingale. It follows that $\int_0^t f dM$ is a martingale for f an elementary function. It is also a continuous martingale because each term is continuous. By Theorem 30.5.3, the maximal estimate and the theorem about the quadratic variation,

$$E \left(\sup_{t \in [0, T]} \left\| \int_0^t f dM \right\|^2 \right) \leq 2E \left(\left\| \int_0^T f dM \right\|^2 \right)$$

Now from the definition of the integral given above, $E \left(\left\| \int_0^T f dM \right\|^2 \right) =$

$$E \left(\sum_{r=0}^{m_n-1} f_r(M(t_{r+1}) - M(t_r)), \sum_{r=0}^{m_n-1} f_r(M(t_{r+1}) - M(t_r)) \right)_H \quad (33.17)$$

because $T \geq t_r$ for each t_r in the partition of the interval. Consider a mixed term in the above product in which $j < k$

$$\begin{aligned} &E(f_k(M(t_{k+1}) - M(t_k)), f_j(M(t_j) - M(t_j))) \\ &= E(E[f_k(M(t_{k+1}) - M(t_k)), f_j(M(t_j) - M(t_j)) | \mathcal{F}_{t_k}]) \\ &= E(E[(M(t_{k+1}) - M(t_k)), f_k^* f_j(M(t_j) - M(t_j)) | \mathcal{F}_{t_k}]) \\ &= E(f_k^* f_j(M(t_j) - M(t_j)), E[(M(t_{k+1}) - M(t_k)) | \mathcal{F}_{t_k}]) = 0 \end{aligned}$$

Now from Lemma 22.5.15 on Hilbert Schmidt operators,

$$\|f_r M\|^2 \leq \|f_r\|^2 \|M\|^2 \leq \|f_r\|_{\mathcal{L}_2}^2 \|M\|^2.$$

Therefore, from 33.17,

$$\begin{aligned} &E \left(\sum_{r=0}^{m_n-1} f_r(M(t_{r+1}) - M(t_r)), \sum_{r=0}^{m_n-1} f_r(M(t_{r+1}) - M(t_r)) \right)_H = \\ &\sum_{r=0}^{m_n-1} E \left(\|f_r(M(t_{r+1}) - M(t_r))\|^2 \right) \leq \sum_{r=0}^{m_n-1} E \left(\|f_r\|^2 \|M(t_{r+1}) - M(t_r)\|^2 \right) \\ &= \sum_{r=0}^{m_n-1} E \left(\|f_r\|_{\mathcal{L}_2}^2 \|M(t_{r+1})\|^2 \right) \\ &+ \sum_{r=0}^{m_n-1} E \left(\|f_r\|_{\mathcal{L}_2}^2 \|M(t_r)\|^2 \right) - 2 \sum_{r=0}^{m_n-1} E \left(\|f_r\|_{\mathcal{L}_2}^2 (M(t_{r+1}), M(t_r)) \right) \quad (33.18) \end{aligned}$$

Consider the mixed term on the end.

$$E \left(\|f_r\|_{\mathcal{L}_2}^2 (M(t_{r+1}), M(t_r)) \right) = E \left(E \left(\|f_r\|_{\mathcal{L}_2}^2 (M(t_{r+1}), M(t_r)) | \mathcal{F}_{t_r} \right) \right)$$

$$= E \left(\|f_r\|_{\mathcal{L}_2}^2 M(t_r), E(M(t_{r+1}) | \mathcal{F}_{t_r}) \right) = E \left(\|f_r\|_{\mathcal{L}_2}^2 \|M(t_r)\|^2 \right)$$

Thus 33.18 reduces to

$$\begin{aligned} & \sum_{r=0}^{m_n-1} E \left(\|f_r\|_{\mathcal{L}_2}^2 \|M(t_{r+1})\|^2 \right) - E \left(\|f_r\|_{\mathcal{L}_2}^2 \|M(t_r)\|^2 \right) \\ &= \int_{\Omega} \sum_{r=0}^{m_n-1} \|f_r\|_{\mathcal{L}_2}^2 \left(\|M(t_{r+1})\|^2 - \|M(t_r)\|^2 \right) dP \\ &= \int_{\Omega} \sum_{r=0}^{m_n-1} \|f_r\|_{\mathcal{L}_2}^2 ([M(t_{r+1})] - [M(t_r)]) dP \end{aligned}$$

because the martingales from the quadratic variation have expectation 0. This has proved the following theorem.

Theorem 33.2.2 *Let f be an elementary function corresponding to the partition $P = \{t_k\}_{k=0}^{m_n}$. Then*

$$\begin{aligned} \frac{1}{2} E \left(\sup_{t \in [0, T]} \left\| \int_0^t f dM \right\|^2 \right) &\leq \int_{\Omega} \sum_{r=0}^{m_n-1} \|f_r\|_{\mathcal{L}_2}^2 ([M(t_{r+1})] - [M(t_r)]) dP \\ &= \int_{\Omega} \int_0^T \|f\|_{\mathcal{L}_2}^2 d[M](t) dP \end{aligned} \quad (33.19)$$

That integral on the right end is just the conventional Stieltjes integral of a step function. Recall that $[M]$ is increasing and continuous.

Theorem 33.2.3 *Let f be uniformly bounded, $\|f(t, \omega)\|_{\mathcal{L}_2(U, H)} \leq K$, continuous in t , and adapted. Also let $M(t)$ be a bounded continuous martingale with values in U and suppose $[M](T) \in L^1(\Omega, P)$. Let f_n be an elementary function approximating f*

$$f_n(t) \equiv \sum_{k=0}^{m_n-1} f(t_k^n) \mathcal{X}_{(t_k^n, t_{k+1}^n]}(t)$$

where $P_n = \{t_k^n\}_{k=0}^{m_n}$ is a partition of $[0, T]$. Assume

$$\|P_n\| \equiv \max \{ |t_{k+1}^n - t_k^n| : k \leq m_n \}.$$

Then there exists a unique continuous bounded martingale denoted as $\int_0^t f dM$ which satisfies

$$\lim_{\|P_n\| \rightarrow 0} \int_0^t f_n dM = \int_0^t f dM \text{ in } \mathcal{M}_T^2(H)$$

where this means: For every $\varepsilon > 0$ there is $\delta > 0$ such that if P_n is a partition having $\|P_n\| < \delta$, then

$$\left\| \int_0^{(\cdot)} f_n dM - \int_0^{(\cdot)} f dM \right\|_{\mathcal{M}_T^2(H)} < \varepsilon.$$

For any such sequence of partitions and approximating elementary functions, there is a set of measure zero N such that if $\omega \notin N$, then $\int_0^t f_n dM(\omega) \rightarrow \int_0^t f dM(\omega)$ uniformly in $t \in [0, T]$. Also, for such f just described,

$$\frac{1}{2} E \left(\sup_{t \in [0, T]} \left\| \int_0^t f dM \right\|^2 \right) \leq \int_{\Omega} \int_0^T \|f\|_{\mathcal{L}_2}^2 d[M](t) dP \quad (33.20)$$

Proof: Letting f_n be an approximating elementary function for f corresponding to P_n a partition with $\|P_n\| \rightarrow 0$, it follows that for each ω , $f_n(t) \rightarrow f(t)$ uniformly in t . From 33.19,

$$\frac{1}{2} \left\| \int_0^t f_n dM(\omega) - \int_0^t f_m dM(\omega) \right\|_{\mathcal{M}_T^2(H)}^2 \leq \int_{\Omega} \int_0^T \|f_n - f_m\|_{\mathcal{L}_2}^2 d[M](t) dP \quad (33.21)$$

By the uniform convergence to f , the inside Stieltjes integral converges to 0 as $n, m \rightarrow \infty$. The integrand is bounded by $4K^2[M](T) \in L^1(\Omega)$ and so one can apply the dominated convergence theorem to conclude that $\{\int_0^t f_n dM(\omega)\}$ is a Cauchy sequence in $\mathcal{M}_T^2(H)$. Therefore, there is a unique $\int_0^t f dM \in \mathcal{M}_T^2(H)$ to which these $\int_0^t f_n dM$ converge. By Proposition 31.7.2 there exists a subsequence converging uniformly in t for all ω off some set of measure zero. The above inequality in 33.21 also implies that $\int_0^{(\cdot)} f dM$ is independent of approximating sequence of elementary functions.

Now $\frac{1}{2} E \left(\sup_{t \in [0, T]} \left\| \int_0^t f_n dM \right\|^2 \right) \leq \int_{\Omega} \int_0^T \|f_n\|_{\mathcal{L}_2}^2 d[M](t) dP$ and so, passing to a limit using the dominated convergence theorem yields 33.20. ■

Definition 33.2.4 Let f be continuous in t , adapted, and uniformly bounded, and let M be a continuous bounded martingale. Then one can define for ω off a set of measure 0, $\int_0^t f dM$ where

$$E \left(\sup_{t \in [0, T]} \left\| \int_0^t f dM - \int_0^t f_n dM \right\|^2 \right) = 0 \quad (33.22)$$

where $\{f_n\}$ is any sequence of elementary functions converging to f .

If f is of bounded variation in addition to being continuous, then the convergence for each ω follows directly from considerations involving Stieltjes integrals without any probabilistic considerations. However, here f is only required to be continuous. Later, this will be relaxed further.

Note that there was no need for f to have values in $\mathcal{L}_2(U, H)$. It would suffice to have $f \in \mathcal{L}(U, H)$ and all of the above would work the same way. You would simply replace $\|f\|_{\mathcal{L}_2(U, H)}$ with $\|f\|_{\mathcal{L}(U, H)}$. The reason for the specialization involves some technical considerations relative to the Wiener process and the need for compactness. Therefore, I have chosen to present it from the beginning with this more specialized case.

Now here is a fundamental lemma about integrals of these elementary functions having to do with stopping times.

Lemma 33.2.5 Suppose M is a continuous bounded martingale having values in a separable Hilbert space U and let f be an elementary function having values in $\mathcal{L}_2(U, H)$ for H a separable Hilbert space. Then the above inequality 33.20 is valid. Also, if σ is a stopping time, then

$$\int_0^t f dM^\sigma = \int_0^{t \wedge \sigma} f dM \quad (33.23)$$

Also,

$$\frac{1}{2} E \left(\sup_{t \in [0, T]} \left\| \int_0^{t \wedge \sigma} f dM \right\|^2 \right) \leq \int_{\Omega} \int_0^T \|f \mathcal{X}_{[0, \sigma]}\|_{\mathcal{L}_2}^2 d[M] dP \quad (33.24)$$

Proof: It only remains to verify 33.23. The first equal sign is obvious from the definition of $\int_0^t f dM$. Both sides equal

$$\sum_{r=0}^{m_n-1} f_r (M(t \wedge \sigma \wedge t_{r+1}) - M(t \wedge \sigma \wedge t_r))$$

Now consider 33.24. From 33.23, the left side is

$$\begin{aligned} \frac{1}{2} E \left(\sup_{t \in [0, T]} \left\| \int_0^t f dM^\sigma \right\|^2 \right) &\leq \int_\Omega \int_0^T \|f\|_{\mathcal{L}_2}^2 d[M^\sigma] dP \\ &= \int_\Omega \int_0^T \|f\|_{\mathcal{L}_2}^2 d[M]^\sigma dP \\ &= \int_\Omega \int_0^T \mathcal{X}_{[0, \sigma]} \|f\|_{\mathcal{L}_2}^2 d[M] dP = \int_\Omega \int_0^T \|\mathcal{X}_{[0, \sigma]} f\|_{\mathcal{L}_2}^2 d[M] dP \end{aligned}$$

This is because when $t > \sigma$, $[M]^\sigma(t) = [M](\sigma)$, a constant. Thus the contribution to the conventional integral is 0 from then on. On the other hand, $[M]^\sigma(t) = [M](t)$ for $t \leq \sigma$ and so the last equation follows. ■

Now here is a nice proposition which is a summary of what has just been discussed along with some other observations.

Proposition 33.2.6 *For $\|f(t, \omega)\|$ bounded by K and continuous in t for each ω having values in $\mathcal{L}_2(U, H)$, and for M a continuous bounded martingale with values in U , $t \rightarrow \int_0^t f dM$ is a continuous martingale with values in H . Assume that $[M](T) \in L^1(\Omega)$. Also the fundamental inequality*

$$\frac{1}{2} E \left(\sup_{t \in [0, T]} \left\| \int_0^t f dM \right\|^2 \right) \leq \int_\Omega \int_0^T \|f\|_{\mathcal{L}_2}^2 d[M] dP \quad (33.25)$$

is valid for f . In addition, $f \rightarrow \int_0^t f dM$ is linear on the linear space of bounded continuous in t adapted functions. If σ is a stopping time,

$$\int_0^{t \wedge \sigma} f dM = \int_0^t f dM^\sigma \quad (33.26)$$

Proof: Since the integral is a limit of integrals of elementary functions for ω off a set of measure zero and since this integral is linear on these functions, the integral is linear. Finally, consider the claim 33.26. From Proposition 31.7.2, and letting $\{f_n\}$ be elementary functions approximating f as above, then by that proposition again, there is a further subsequence still denoted with n such that for a.e. ω ,

$$\int_0^t f dM^\sigma = \lim_{n \rightarrow \infty} \int_0^t f_n dM^\sigma = \lim_{n \rightarrow \infty} \int_0^{t \wedge \sigma} f_n dM = \int_0^{t \wedge \sigma} f dM \quad \blacksquare$$

Note how this does not require f to be of bounded variation. If f were of bounded variation, you would get pointwise convergence of the Stieltjes sums for $\int_0^t f_n dM$ to $\int_0^t f dM$ as a consequence of simple considerations involving Stieltjes integrals. Then the new information is that the Stieltjes integral $\int_0^t f dM$ for f of bounded variation is a continuous martingale.

33.3 More on Stopping Times

Next I want to consider $\int_0^t f \mathcal{X}_{[0,\sigma]} dM$ in the case of elementary functions for f a continuous in t , adapted, and bounded. This involves the same process as earlier and it is analogous to the traditional definition of the Stieltjes integral. First we approximate with an elementary function and then pass to a limit.

Proposition 33.3.1 *Let σ be a stopping time and let $\{t_k\}_{k=0}^{m_n}$ be partition points of $[0, T]$. Also define the discrete approximation of σ*

$$\sigma_n(\omega) \equiv \sum_{k=0}^{m_n-1} t_k \mathcal{X}_{\sigma^{-1}((t_k, t_{k+1}])}(\omega)$$

Then for f an elementary function with respect to $\{t_k\}_{k=0}^n$ points as described above, it follows that $f \mathcal{X}_{[0,\sigma_n]}$ is an elementary function and σ_n is a stopping time. Also

$$\int_0^t f \mathcal{X}_{[0,\sigma_n]} dM = \int_0^{t \wedge \sigma_n} f dM$$

Proof: First, why is σ_n a stopping time? Consider $[\sigma_n \leq t]$. Say $t \in (t_k, t_{k+1}]$. In case, $t = t_{k+1}$, $[\sigma_n \leq t] = [\sigma \leq t] \in \mathcal{F}_t$. Otherwise, $[\sigma_n \leq t] = [\sigma \leq t_k] \in \mathcal{F}_{t_k} \subseteq \mathcal{F}_t$. Thus σ_n is indeed a stopping time. Now

$$f \mathcal{X}_{[0,\sigma_n]}(t) = \sum_{k=0}^{m_n-1} f_k \mathcal{X}_{[0,\sigma_n]}(t) \mathcal{X}_{(t_k, t_{k+1}]}(t).$$

t is somewhere. Say $t \in (t_k, t_{k+1}]$. Then consider $f_k \mathcal{X}_{[0,\sigma_n]}(t)$. For this t , this term is nonzero if and only if $\omega \in [t \leq \sigma_n]$ if and only if $\omega \in [t_k < \sigma_n] \in \mathcal{F}_{t_k}$. Thus this term is the indicator function of a set in \mathcal{F}_{t_k} for all $t \in (t_k, t_{k+1}]$. It follows $f \mathcal{X}_{[0,\sigma_n]}$ is an elementary function and can be written as $\sum_{k=0}^{m_n-1} f_k \mathcal{X}_{[t_k < \sigma_n]} \mathcal{X}_{(t_k, t_{k+1}]}(t)$, so its integral is

$$\begin{aligned} & \sum_{k=0}^{m_n-1} f_k \mathcal{X}_{[t_k < \sigma_n]} (M(t \wedge t_{k+1}) - M(t \wedge t_k)) \\ &= \sum_{k=0}^{m_n-1} f_k (M(t \wedge t_{k+1} \wedge \sigma_n) - M(t \wedge t_k \wedge \sigma_n)) \end{aligned}$$

because if $\sigma_n > t_k$ in the k^{th} term and the term is nonzero, then $\sigma_n = t_{k+1}$. If $\sigma_n \leq t_k$, the k^{th} term on the left is 0 and on the right that term is

$$f_k (M(t \wedge \sigma_n) - M(t \wedge \sigma_n)),$$

also zero. Now the right side in the above is just $\int_0^{t \wedge \sigma_n} f dM$ and the left side is defined earlier as $\int_0^t f \mathcal{X}_{[0,\sigma_n]} dM$. ■

Lemma 33.3.2 *Let f be bounded, adapted, and continuous in t and let σ be a stopping time with finite values in $[0, T]$. Also assume $[M](T) \in L^1(\Omega, P)$. Then letting σ_n be as above,*

$$\sigma_n(\omega) \equiv \sum_{k=0}^{m_n-1} t_k \mathcal{X}_{\sigma^{-1}((t_k, t_{k+1}])}(\omega),$$

$f_n \mathcal{X}_{[0, \sigma_n]} \rightarrow f \mathcal{X}_{[0, \sigma]}$ in $L^2(\Omega \times [0, T]; \mathcal{L}_2(U, H))$. Here $f_n(t)$ is the elementary function $\sum_{k=0}^{m_n-1} f(t_k^n) \mathcal{X}_{(t_k^n, t_{k+1}^n]}(t)$ where $\{t_k^n\}$ is a partition P_n with $\|P_n\| \rightarrow 0$ so that f_n converges uniformly to f for each ω .

Proof: Note that $|\sigma_n(\omega) - \sigma(\omega)| \leq \|P_n\|$. Then

$$\begin{aligned} \|f_n \mathcal{X}_{[0, \sigma_n]} - f \mathcal{X}_{[0, \sigma]}\| &\leq \|(f_n - f) \mathcal{X}_{[0, \sigma_n]}\| + \|f \mathcal{X}_{[\sigma_n, \sigma]}\| \\ \|f_n \mathcal{X}_{[0, \sigma_n]} - f \mathcal{X}_{[0, \sigma]}\|^2 &\leq 2 \left(\|(f_n - f) \mathcal{X}_{[0, \sigma_n]}\|^2 + \|f \mathcal{X}_{[\sigma_n, \sigma]}\|^2 \right) \end{aligned}$$

Thus

$$\begin{aligned} &\int_{\Omega} \int_0^T \|f_n \mathcal{X}_{[0, \sigma_n]} - f \mathcal{X}_{[0, \sigma]}\|^2 d[M] dP \\ &\leq 2 \int_{\Omega} \int_0^T \|f_n - f\|^2 d[M] dP + 2K \int_{\Omega} \int_0^T \mathcal{X}_{[\sigma_n, \sigma]} d[M] dP \\ &\leq 2 \int_{\Omega} \int_0^T \|f_n - f\|^2 d[M] dP + 2K \int_{\Omega} [M](\sigma) - [M](\sigma_n) dP \quad (33.27) \end{aligned}$$

the integrands $\int_0^T \|f_n - f\|^2 d[M]$ and $[M](\sigma) - [M](\sigma_n)$ both converge to 0 a.e. ω as $n \rightarrow \infty$ thanks to continuity of $[M]$. Also the assumption that $[M](T)$ is in L^1 along with the boundedness of f imply these integrands are uniformly integrable. Hence we can use the Vitali convergence theorem and conclude that the limit of 33.27 is 0. ■

Now note that if you fix ω , $f_n \mathcal{X}_{[0, \sigma_n]}(t) \rightarrow f \mathcal{X}_{[0, \sigma]}(t)$ for each $t < \sigma(\omega)$. Thus, if f is of bounded variation as well as being continuous in each t , standard Stieltjes integral considerations involving continuity of f and M show that $\int_0^t f_n \mathcal{X}_{[0, \sigma_n]}(t) dM \rightarrow \int_0^t f \mathcal{X}_{[0, \sigma]}(t) dM$ with no probabilistic complications at all.

Definition 33.3.3 Let f be adapted, continuous in t , and bounded. Let M be a continuous bounded martingale. Also let σ be a stopping time and let σ_n be the discreet approximation above relative to partitions $P_n = \{t_k^n\}_{k=0}^{m_n-1}$ where $\|P_n\| \rightarrow 0$ and $\{f_n\}$ the sequence of elementary functions approximating f ,

$$f_n(t) \equiv \sum_{k=0}^{m_n-1} f(t_k^n) \mathcal{X}_{(t_k^n, t_{k+1}^n]}(t),$$

Then there exists a martingale denoted as $\int_0^t f \mathcal{X}_{[0, \sigma]} dM$ such that

$$\int_0^t f \mathcal{X}_{[0, \sigma]} dM = \lim_{n \rightarrow \infty} \int_0^t f_n \mathcal{X}_{[0, \sigma_n]} dM \text{ in } M_T^2(H)$$

This is something new. Earlier we had $\int_0^t f dM$ defined where f is continuous on $[0, T]$. We also have $\int_0^{t \wedge \sigma} f dM$ defined where σ is a stopping time and f is continuous on $[0, T]$. However, $f \mathcal{X}_{[0, \sigma]}$ is not necessarily continuous on $[0, T]$. It is continuous on $[0, T \wedge \sigma]$ so the time intervals are changing as a function of ω .

To begin with, we can stop the martingale $\int_0^t f_n dM$ with the stopping time σ , this denoted as $\int_0^{t \wedge \sigma} f_n dM$. Then

$$\begin{aligned} &\int_0^{t \wedge \sigma} f_n dM - \int_0^{t \wedge \sigma_n} f_n dM \\ &= \sum_{k=0}^{m_n-1} f(t_k^n) \left(\begin{aligned} &(M(t \wedge \sigma \wedge t_{k+1}) - M(t \wedge \sigma \wedge t_k)) \\ &- (M(t \wedge \sigma_n \wedge t_{k+1}) - M(t \wedge \sigma_n \wedge t_k)) \end{aligned} \right) \end{aligned}$$

$$= \sum_{k=0}^{m_n-1} f(t_k^n) \begin{pmatrix} (M(t \wedge \sigma \wedge t_{k+1}) - M(t \wedge \sigma_n \wedge t_{k+1})) \\ - (M(t \wedge \sigma \wedge t_k) - (M(t \wedge \sigma_n \wedge t_k))) \end{pmatrix}$$

Now from maximal theorems,

$$\begin{aligned} & E \left(\sup_{t \in [0, T]} \left\| \int_0^{t \wedge \sigma} f_n dM - \int_0^{t \wedge \sigma_n} f_n dM \right\|^2 \right) \\ & \leq E \left(\left\| \int_0^{T \wedge \sigma} f_n dM - \int_0^{T \wedge \sigma_n} f_n dM \right\|^2 \right) \\ & = E \left(\sum_{k=0}^{m_n-1} \left\| f(t_k^n) (M(\sigma \wedge t_{k+1}) - M(\sigma_n \wedge t_{k+1})) \right\|^2 \right. \\ & \quad \left. + \left\| f(t_k^n) (M(\sigma \wedge t_k) - (M(\sigma_n \wedge t_k))) \right\|^2 \right) + ?? \end{aligned} \quad (33.28)$$

where ?? is -2 times the expectation of a sum of mixed terms which can be written in the following form after noticing that $M(\sigma \wedge t_k) - (M(\sigma_n \wedge t_k))$ is \mathcal{F}_{t_k} measurable.

$$\begin{aligned} & E \left(\begin{matrix} f(t_k^n)^* f(t_k^n) M(\sigma \wedge t_k) \\ - (M(\sigma_n \wedge t_k)), E(M(\sigma \wedge t_{k+1}) - M(\sigma_n \wedge t_{k+1}) | \mathcal{F}_{t_k}) \end{matrix} \right) \\ & = E(f(t_k^n)^* f(t_k^n) M(\sigma \wedge t_k) - (M(\sigma_n \wedge t_k)), M(\sigma \wedge t_k) - M(\sigma_n \wedge t_k)) \\ & = E(\|f(t_k^n) (M(\sigma \wedge t_k) - M(\sigma_n \wedge t_k))\|^2) \end{aligned}$$

Therefore 33.28 reduces to

$$\begin{aligned} & E \left(\sum_{k=0}^{m_n-1} \left\| f(t_k^n) (M(\sigma \wedge t_{k+1}) - M(\sigma_n \wedge t_{k+1})) \right\|^2 \right. \\ & \quad \left. - \left\| f(t_k^n) (M(\sigma \wedge t_k) - (M(\sigma_n \wedge t_k))) \right\|^2 \right) \\ & = E(\|f(T) (M(T \wedge \sigma) - M(T \wedge \sigma_n))\|^2) \end{aligned}$$

This converges to 0 because the integrand converges to 0 since $\sigma_n(\omega) \rightarrow \sigma(\omega)$ and the integrand is uniformly bounded by assumption. Thus this has shown that

$$E \left(\sup_{t \in [0, T]} \left\| \int_0^{t \wedge \sigma} f_n dM - \int_0^{t \wedge \sigma_n} f_n dM \right\|^2 \right) \rightarrow 0.$$

This has shown the following technical lemma.

Lemma 33.3.4 *Letting f be bounded and continuous in t and adapted and letting M be a bounded continuous martingale, and σ_n the discrete approximation of a stopping time σ and f_n the elementary function approximating f as described above, then it follows that*

$$E \left(\sup_{t \in [0, T]} \left\| \int_0^{t \wedge \sigma} f_n dM - \int_0^{t \wedge \sigma_n} f_n dM \right\|^2 \right) \rightarrow 0.$$

Theorem 33.3.5 *Definition 33.3.3, is well defined and $\int_0^t f \mathcal{X}_{[0, \sigma]} dM$ is a martingale equal to $\int_0^{t \wedge \sigma} f dM$. Also, there is a set of measure zero N and a subsequence, still denoted as n such that for $\omega \notin N$,*

$$\int_0^t f_n \mathcal{X}_{[0, \sigma_n]} dM(\omega) \rightarrow \int_0^t f \mathcal{X}_{[0, \sigma]} dM(\omega) \quad (33.29)$$

uniformly in $t \in [0, T]$.

Proof: From Proposition 31.7.2, $\int_0^{t \wedge \sigma} f_n dM \rightarrow \int_0^{t \wedge \sigma} f dM$ in $\mathcal{M}_T^2(H)$. From Lemma 33.3.2 it follows that also

$$\int_0^{t \wedge \sigma_n} f_n dM = \int_0^t f_n \mathcal{X}_{[0, \sigma_n]} dM \rightarrow \int_0^{t \wedge \sigma} f dM$$

in $\mathcal{M}_T^2(H)$. Therefore, $\{\int_0^t f_n \mathcal{X}_{[0, \sigma_n]} dM\}$ is a Cauchy sequence in $\mathcal{M}_T^2(H)$ and so converges to a continuous martingale denoted as $\int_0^t f \mathcal{X}_{[0, \sigma]} dM$. Any two sequences of approximating elementary functions lead to the same $\int_0^t f \mathcal{X}_{[0, \sigma]} dM$. Thanks to Lemma 33.3.2 and the inequality

$$\left\| \int_0^{t \wedge \sigma} f_n dM - \int_0^{t \wedge \sigma} \hat{f}_n dM \right\|_{\mathcal{M}_T^2(H)}^2 \leq 2 \int_{\Omega} \int_0^T \|f_n - \hat{f}_n\|_{\mathcal{L}_2}^2 d[M] dP$$

in which the right side converges to 0. By Proposition 31.7.2 again, there is a set of measure zero N such that if $\omega \notin N$, then 33.29 holds. ■

Again, if f is not just continuous and adapted but is also of bounded variation, the convergence for each ω follows from Stieltjes integration theory. Essentially, what the above shows is that even in this case, the integral is a continuous martingale. Also, the above gives meaning to the expression $\int_0^t f \mathcal{X}_{[0, \sigma]} dM(\omega)$ as a limit in $\mathcal{M}_T^2(H)$ of integrals of appropriate elementary functions just as $\int_0^t f dM(\omega)$ was a limit of integrals of appropriate elementary functions. It is the same thing you see with Stieltjes integration but here it is much more elaborate. In place of bounded variation you have adapted and the integrator function is now a martingale.

33.4 Local Martingales as Integrators

Now suppose M is a local martingale. This means that there is a localizing sequence of stopping times $\{\tau_n\}$, $\tau_n \rightarrow \infty$, such that M^{τ_n} is a martingale. Recall Proposition 32.3.3 which says that we can always assume the localizing sequence τ_n makes M^{τ_n} uniformly bounded as well as a martingale.

I will use this fact whenever convenient from now on. We know that $\int_0^{t \wedge \tau} f dM \equiv \int_0^t f dM^{\tau}$ if M is a bounded martingale so the following definition will simply extend this idea to give a definition for the stochastic integral in the case where M is continuous but maybe not bounded.

Definition 33.4.1 Let M be a continuous local martingale with a localizing sequence τ_n for which M^{τ_n} is a bounded martingale. Let f be bounded and continuous in t and adapted. Then $\int_0^t f dM$ is defined as follows.

$$\int_0^{t \wedge \tau_n} f dM \equiv \int_0^t f dM^{\tau_n}$$

thus for a given ω , $\int_0^t f dM = \lim_{n \rightarrow \infty} \int_0^t f dM^{\tau_n}$.

Lemma 33.4.2 The above definition does define an adapted process $\int_0^t f dM$ which is a continuous local martingale.

Proof: First note that if σ_n is another localizing sequence, $\int_0^t f dM(\omega)$ is the same when defined from either localizing sequence. The reason is that for either sequence, there is n

such that if $m \geq n$, then for a given ω , $\tau_m(\omega), \sigma_m(\omega)$ are both ∞ . It follows upon using approximation with elementary functions and passing to a limit using a suitable subsequence of elementary functions that for such τ_m, σ_m ,

$$\int_0^{t \wedge \tau_m} f dM(\omega), \int_0^{t \wedge \sigma_m} f dM(\omega) = \lim_{p \rightarrow \infty} \sum_{k=0}^{p-1} f(t_k^p) (M(t \wedge t_{k+1}^p) - M(t \wedge t_k^p))(\omega)$$

this limit being independent of the localizing sequence used. The reason it is a local continuous martingale is that M^{τ_n} is a bounded continuous martingale and so

$$\left(\int_0^t f dM \right)^{\tau_n} \equiv \int_0^{t \wedge \tau_n} f dM \equiv \int_0^t f dM^{\tau_n}$$

is a martingale. ■

The idea is that if you know it at $t \wedge \tau_n$ for all τ_n where $\tau_n \rightarrow \infty$, then you know it at t because you can simply pick τ_n larger than t and from the above, it doesn't matter which localizing sequence you use.

Now suppose M is just a local martingale so there is a localizing sequence of stopping times $\{\tau_n\}$ such that M^{τ_n} is a bounded martingale and suppose $[M](T) \in L^1(\Omega)$. Let f_n be elementary functions converging to f a bounded, adapted, continuous in t function, convergence uniform in t for each ω . Then from the fundamental inequality above and using the formulas for stopping times,

$$\begin{aligned} \frac{1}{2} E \left(\sup_{t \in [0, T]} \left\| \int_0^{t \wedge \tau_p} f_n dM \right\|^2 \right) &= \frac{1}{2} E \left(\sup_{t \in [0, T]} \left\| \int_0^t f_n dM^{\tau_p} \right\|^2 \right) \\ &\leq \int_{\Omega} \int_0^T \|f_n\|^2 d[M]^{\tau_p}(t) dP \leq \int_{\Omega} \int_0^T \|f_n\|^2 d[M](t) dP \end{aligned}$$

Since $\|f\|$ is assumed bounded and $[M](T) \in L^1$, it follows that there is a dominating function on the right, namely $K^2 [M](T)$ where $K \geq \|f(t, \omega)\|$ for all (t, ω) . Let $n \rightarrow \infty$ and use the dominated convergence on the right and either Fatou's lemma or monotone convergence theorem on the left to obtain

$$\frac{1}{2} E \left(\sup_{t \in [0, T]} \left\| \int_0^t f dM \right\|^2 \right) \leq \int_{\Omega} \int_0^T \|f\|^2 d[M](t) dP$$

The reason the monotone convergence theorem applies is that on the left, the stopping time effectively restricts the values of t over which the sup is taken until $\tau_p \geq t$. You could also apply Fatou's lemma.

This has shown the following proposition.

Proposition 33.4.3 *Let f be adapted, continuous in t and bounded. Also let M be a local martingale with $[M](T) \in L^1(\Omega)$. Then*

$$\frac{1}{2} E \left(\sup_{t \in [0, T]} \left\| \int_0^t f dM \right\|^2 \right) \leq \int_{\Omega} \int_0^T \|f\|^2 d[M](t) dP$$

If f is bounded, continuous, and of bounded variation, and adapted and M is a local martingale, this shows that the Stieltjes integral $\int_0^t f dM$ is a local martingale.

It is clear that if M is a local martingale with localizing sequence τ_p and if σ is a stopping time, then M^σ is also a local martingale with localizing sequence τ_p because $(M^\sigma)^{\tau_p} = (M^{\tau_p})^\sigma$, the latter being a bounded martingale. Now for f continuous in t and bounded and adapted,

$$\int_0^{t \wedge \sigma \wedge \tau_p} f dM \equiv \int_0^t f dM^{\sigma \wedge \tau_p} \equiv \int_0^{t \wedge \tau_p} f dM^\sigma$$

Therefore, from the definition, whenever M is a local martingale and f bounded and continuous in t and adapted,

$$\int_0^{t \wedge \sigma} f dM = \int_0^t f dM^\sigma \quad (33.30)$$

33.5 The Stochastic Integral and the Quadratic Variation

In this simple case of the above, you have a bounded martingale M with values in U and you have $f \in U'$. Thus $f \in \mathcal{L}(U, \mathbb{R})$. Is f actually in $\mathcal{L}_2(U, \mathbb{R})$? By Riesz representation theorem, there is $x \in U$ such that $Rx = f$. Then if $\{g_k\}$ is an orthonormal basis for U , $\sum_k |f(g_k)|^2 = \sum_k |(g_k, Rx)|^2 = \sum_k |(g_k, x)|^2 = \|x\|^2$ because it is just the sum of the squares of the Fourier coefficients of x . Thus $f \in \mathcal{L}_2(U, \mathbb{R})$.

Now an example of a continuous in t , adapted, bounded function in $\mathcal{L}_2(U, \mathbb{R})$ is just $RM(t) \equiv f(t)$. Therefore, it makes perfect sense to consider $\int_0^t (RM) dM$. Let $\|P_n\| \rightarrow 0$ and let the stochastic integral of an elementary function $f_n(t) = \sum_{k=0}^{m_n-1} RM(t_k^n) \mathcal{X}_{(t_k^n, t_{k+1}^n]}(t)$ be of the form

$$\sum_{k=0}^{m_n-1} RM(t_k^n) (M(t \wedge t_{k+1}^n) - M(t \wedge t_k^n))$$

where $\{t_k^n\}_{k=0}^{m_n}$ is this partition P_n .

ALWAYS assume in this that the partitions are nested, $P_n \subseteq P_{n+1}$.

$$\begin{aligned} Q_n(t) &\equiv \sum_{k=0}^{m_n-1} \|M(t \wedge t_{k+1}^n) - M(t \wedge t_k^n)\|_U^2 \\ &= \sum_{k=0}^{m_n-1} \|M(t \wedge t_{k+1}^n)\|^2 + \|M(t \wedge t_k^n)\|^2 - 2(M(t \wedge t_{k+1}^n), M(t \wedge t_k^n)) \\ &= \sum_{k=0}^{m_n-1} \|M(t \wedge t_{k+1}^n)\|^2 - \|M(t \wedge t_k^n)\|^2 - 2(M(t \wedge t_k^n), M(t \wedge t_{k+1}^n) - M(t \wedge t_k^n)) \\ &= \|M(t)\|^2 - 2 \int_0^t (RM_n) dM \end{aligned}$$

Then passing to a limit, then $Q_n(t) \rightarrow Q(t)$ in $L^2(\Omega)$ because $2 \int_0^t (RM_n) dM$ converges in $\mathcal{M}_T^2(\mathbb{R})$. Using a subsequence, we can also get uniform convergence in t for all ω off a set of measure zero. Thus Q is increasing. It follows

$$Q(t) = \|M(t)\|^2 - 2 \int_0^t (RM) dM = [M](t) + N(t) - 2 \int_0^t (RM) dM$$

and so $Q(t) - [M](t)$ equals a martingale. Thus from Lemma 32.2.1, $Q(t) - [M](t) = 0$. This proves the first part of the following important result.

Theorem 33.5.1 *Let H be a Hilbert space and suppose $(M, \mathcal{F}_t), t \in [0, T]$ is a uniformly bounded continuous martingale with values in H . Also let $\{t_k^n\}_{k=1}^{m_n}$ be a sequence of partitions satisfying*

$$\lim_{n \rightarrow \infty} \max \{ |t_i^n - t_{i+1}^n|, i = 0, \dots, m_n \} = 0, \{t_k^n\}_{k=1}^{m_n} \subseteq \{t_k^{n+1}\}_{k=1}^{m_{n+1}}.$$

Then

$$[M](t) = \lim_{n \rightarrow \infty} \sum_{k=0}^{m_n-1} \|M(t \wedge t_{k+1}^n) - M(t \wedge t_k^n)\|_H^2$$

the limit taking place in $L^2(\Omega)$. In case M is just a continuous local martingale, the above limit happens in probability.

Proof: It only remains to show the claim about the case where M is a local martingale. Suppose M is only a continuous local martingale. By Proposition 32.3.3 there exists an increasing localizing sequence $\{\tau_k\}$ such that M^{τ_k} is a uniformly bounded martingale. Then

$$P(\cup_{k=1}^{\infty} [\tau_k = \infty]) = 1$$

As above, let

$$Q_n(t) \equiv \sum_{k=0}^{m_n-1} \|M(t \wedge t_{k+1}^n) - M(t \wedge t_k^n)\|_H^2$$

where there are m_n points in P_n where as before, $P_n \subseteq P_{n+1}$ for all n .

Let $\eta, \varepsilon > 0$ be given. Then there exists k large enough that $P([\tau_k = \infty]) > 1 - \eta/2$. This is because the sets $[\tau_k = \infty]$ increase to Ω other than a set of measure zero. Then for this k ,

$$[|Q_n^{\tau_k} - [M]^{\tau_k}(t)| > \varepsilon] \cap [\tau_k = \infty] = [|Q_n - [M](t)| > \varepsilon] \cap [\tau_k = \infty]$$

Thus

$$\begin{aligned} P(|Q_n - [M](t)| > \varepsilon) &\leq P(|Q_n - [M](t)| > \varepsilon \cap [\tau_k = \infty]) \\ &\quad + P([\tau_k < \infty]) \\ &\leq P(|Q_n^{\tau_k} - [M]^{\tau_k}(t)| > \varepsilon) + \eta/2 \end{aligned}$$

The convergence in probability of $Q_n^{\tau_k}(t)$ to $[M]^{\tau_k}(t)$ follows from the convergence in $L^2(\Omega)$ shown earlier for bounded martingales, and so if n is large enough, the right side of the above inequality is less than $\eta/2 + \eta/2 = \eta$. Since η was arbitrary, this proves convergence in probability. ■

33.6 The Case of $f \in L^2(\Omega \times [0, T]; \mathcal{L}_2(U, H))$

At this point I will discontinue the general treatment in terms of arbitrary martingales and suppose that $[M](t) = F(t)$ a continuous increasing function which depends only on t and not on ω . In fact, the most interest is centered on the Wiener process in which $[M](t) = at$ for $a > 0$ and this is the case considered here. Of course a does not matter so we will simply assume $[M](t) = t$.

It is not necessary to have f be uniformly bounded. Instead, one can consider $f \in L^2(\Omega \times [0, T]; \mathcal{L}_2(U, H))$ where f is progressively measurable. I considered the case where f is continuous in t above because it is a convenient way to tie this in to the ordinary theory of Stieltjes integrals and to point out that these standard objects do deliver martingales in some reasonable cases.

As before, I will first consider the case where M is a bounded martingale and then extend to the case where M is a local martingale. As before, it is all based on the fundamental inequality

$$\frac{1}{2}E\left(\sup_{t \in [0, T]} \left\| \int_0^t f dM \right\|^2\right) \leq E\left(\left\| \int_0^T f dM \right\|^2\right) \leq \int_{\Omega} \int_0^T \|f\|_{\mathcal{L}_2}^2 dt dP \quad (33.31)$$

which was shown to hold for all adapted f continuous in t and also bounded, which implies the integral on the right is finite.

Let $f \in L^2(\Omega \times [0, T]; \mathcal{L}_2(U, H))$ be adapted and let $f(t, \omega)$ be extended as 0 for t off $[0, T]$.

Definition 33.6.1 Let $f_n(t) \equiv n \int_{t-1/n}^t f(s) ds$.

Lemma 33.6.2 f_n is adapted and $t \rightarrow f_n(t)$ is continuous for a.e. ω .

Proof: First consider the claim about continuity. For each ω off a set of measure zero, $f \in L^2([0, T]; \mathcal{L}_2)$ and so, for such ω

$$\begin{aligned} \|f_n(t) - f_n(\hat{t})\|_{\mathcal{L}_2}^2 &= \left\| n \int_{t-1/n}^t f(s) ds - n \int_{\hat{t}-1/n}^{\hat{t}} f(s) ds \right\|_{\mathcal{L}_2}^2 dP \\ &= n \left\| \int_{t-1/n}^{\hat{t}-1/n} f(s) ds + \int_{\hat{t}}^t f(s) ds \right\|_{\mathcal{L}_2}^2 dP \\ &\leq 2n \left(\left\| \int_{t-1/n}^{\hat{t}-1/n} f(s) ds \right\|_{\mathcal{L}_2}^2 + \left\| \int_{\hat{t}}^t f(s) ds \right\|_{\mathcal{L}_2}^2 \right) \leq 8n(t - \hat{t}) \|f\|_{L^2([0, T]; \mathcal{L}_2)}^2 \end{aligned}$$

It follows that $\omega \rightarrow n \int_{t-1/n}^t f(s) ds$ is \mathcal{F}_t measurable because $\mathcal{H}_{[0, t]} f$ is $\mathcal{F}_t \times \mathcal{B}([0, T])$ measurable by assumption that f is progressively measurable. ■

Observe that

$$\begin{aligned} \|f_n - f\|_{L^2(\Omega \times [0, T]; \mathcal{L}_2(U, H))} &\equiv \left(\int_0^T \int_{\Omega} \|f_n - f\|_{\mathcal{L}_2}^2 dP dt \right)^{1/2} \\ &= \left(\int_{\Omega} \int_0^T \left\| n \int_{-1/n}^0 (f(t+s) - f(t)) ds \right\|_{\mathcal{L}_2}^2 dt dP \right)^{1/2} \end{aligned}$$

From Minkowski's inequality,

$$\leq n \int_{-1/n}^0 \left(\int_{\Omega} \int_0^T \|f(t+s) - f(t)\|^2 dt dP \right)^{1/2} ds \quad (33.32)$$

so

$$\|f_n - f\|_{L^2(\Omega \times [0, T]; \mathcal{L}_2(U, H))}^2 \leq \int_{\Omega} n \int_{-1/n}^0 \int_0^T \|f(t+s) - f(t)\|^2 dt ds dP$$

Now

$$\begin{aligned}
& n \int_{-1/n}^0 \int_0^T \|f(t+s) - f(t)\|^2 dt ds \\
& \leq 2n \int_{-1/n}^0 \int_0^T \|f(t+s)\|^2 dt ds + 2n \int_{-1/n}^0 \int_0^T \|f(t)\|^2 dt ds \\
& \leq 4 \int_0^T \|f(t)\|^2 ds
\end{aligned}$$

which by definition is in $L^1(\Omega)$. Therefore, the integrands

$$n \int_{-1/n}^0 \int_0^T \|f(t+s) - f(t)\|^2 dt ds$$

converge to 0 by continuity of translation in L^2 and are uniformly integrable. By Vitali convergence theorem, $\lim_{n \rightarrow \infty} \|f_n - f\|_{L^2(\Omega \times [0, T]; \mathcal{L}_2(U, H))}^2 = 0$.

I want to use the fundamental inequality 33.31 which has only been presented above for f bounded. Therefore, let $g_n(t, \omega) \equiv P_{m_n}(f(t, \omega))$ where P_{m_n} is the projection onto $\overline{B(0, m_n)}$ in the Hilbert space $\mathcal{L}_2(U, H)$. As follows from the definition, one can obtain an inner product for the norm in this Banach space in the form $(f, g) \equiv \sum_{k=1}^{\infty} (f(e_k), g(e_k))_H$ where $\{e_k\}$ is some orthonormal basis for U . Now P_{m_n} is Lipschitz continuous and if m_n is large enough, $\|g_n - f\|_{L^2(\Omega \times [0, T]; \mathcal{L}_2(U, H))} < 2^{-n}$ and so $g_n(t, \omega)$ is bounded and continuous and adapted and

$$\lim_{n \rightarrow \infty} \|g_n - f\|_{L^2(\Omega \times [0, T]; \mathcal{L}_2(U, H))} = 0 \quad (33.33)$$

With this preparation, here is the main result.

Theorem 33.6.3 *Let f be adapted and in $L^2(\Omega \times [0, T]; \mathcal{L}_2(U, H))$. Then there exists a sequence $\{g_n\}$ of adapted functions continuous in t such that*

$$\lim_{n \rightarrow \infty} \|g_n - f\|_{L^2(\Omega \times [0, T]; \mathcal{L}_2(U, H))}^2 = \lim_{n \rightarrow \infty} \int_{\Omega} \int_0^T \|g_n - f\|^2 dt dP = 0$$

Also it follows that $\int_0^t g_n dM$ is a Cauchy sequence in $\mathcal{M}_T^2(H)$ converging to a continuous martingale denoted as $\int_0^t f dM$ in $\mathcal{M}_T^2(H)$. In addition, for $f \in L^2(\Omega \times [0, T]; \mathcal{L}_2(U, H))$, adapted, the fundamental inequality holds.

$$\frac{1}{2} E \left(\sup_{t \in [0, T]} \left\| \int_0^t f dM \right\|^2 \right) \leq \int_{\Omega} \int_0^T \|f\|_{\mathcal{L}_2}^2 dt dP$$

If M is only a local martingale, then the same inequality is valid. Also, if σ is any stopping time,

$$\int_0^{t \wedge \sigma} f dM = \int_0^t f dM^\sigma$$

Proof: It follows from the above argument there exists a sequence of adapted continuous, bounded, functions converging to f in $L^2(\Omega \times [0, T]; \mathcal{L}_2(U, H))$. Therefore,

$$\lim_{m, n \rightarrow \infty} \frac{1}{2} E \left(\sup_{t \in [0, T]} \left\| \int_0^t (g_n - g_m) dM \right\|^2 \right) \leq \lim_{m, n \rightarrow \infty} \int_{\Omega} \int_0^T \|g_n - g_m\|^2 dt dP = 0$$

and by completeness of $\mathcal{M}_T^2(H)$, $\int_0^t g_n dM$ converges to a continuous martingale which I can call $\int_0^t f dM$. It is clear that any two sequences give the same result from the inequality satisfied. Therefore, the stochastic integral $\int_0^t f dM$ is well defined. Also from the theory of $\mathcal{M}_T^2(H)$, there is a subsequence for which $\int_0^t g_n dM$ converges uniformly in t to $\int_0^t f dM$ off some set of measure zero.

In case M is only a local martingale, we see from approximating f with continuous in t and adapted and bounded functions g_n as above that the appropriate way to define $\int_0^t f dM$ is as $\int_0^{t \wedge \sigma_n} f dM \equiv \int_0^t f dM^{\sigma_n}$ where $\{\sigma_n\}$ is a localizing sequence for M .

The quadratic variation of M^{σ_n} is no more than the quadratic variation of M and so

$$\begin{aligned} & \lim_{m, n \rightarrow \infty} \frac{1}{2} E \left(\sup_{t \in [0, T]} \left\| \int_0^{t \wedge \sigma_n} (g_n - g_m) dM \right\|^2 \right) \\ & \leq \lim_{m, n \rightarrow \infty} \int_{\Omega} \int_0^T \|g_n - g_m\|^2 dt dP = 0 \end{aligned}$$

Thus we can obtain $\int_0^t f dM^{\sigma_n}$ as a limit in $\mathcal{M}_T^2(H)$ as just done. Then one can define $\int_0^t f dM \equiv \lim_{n \rightarrow \infty} \int_0^t f dM^{\sigma_n}$ where σ_n is a localizing sequence for M . Also, we can pass to a limit as $n \rightarrow \infty$ using the monotone convergence theorem in the inequality

$$\frac{1}{2} E \left(\sup_{t \in [0, T]} \left\| \int_0^{t \wedge \sigma_n} f dM^{\sigma_n} \right\|^2 \right) \leq \int_{\Omega} \int_0^T \|f\|^2 dt dP$$

The stopping time has the effect of restricting the time interval, so as σ_n increases, one is taking sup over a larger set. That is why the monotone convergence theorem applies on the left side. Then

$$\frac{1}{2} E \left(\sup_{t \in [0, T]} \left\| \int_0^t f dM \right\|^2 \right) \leq \int_{\Omega} \int_0^T \|f\|^2 dt dP$$

As to the last claim about stopping times, it works for f bounded and continuous in t and adapted and M a martingale. Therefore, it also works for

$$f \in L^2(\Omega \times [0, T]; \mathcal{L}_2(U, H)).$$

In general, when M is only a local martingale with localizing sequence stopping times $\{\tau_n\}$, then M^{σ} is also a local martingale with localizing sequence $\{\tau_n\}$. I need to show that

$$\int_0^{t \wedge \sigma} f dM = \int_0^t f dM^{\sigma}$$

as local martingales. I need to show that

$$\int_0^{t \wedge \sigma \wedge \tau_n} f dM = \int_0^t f d(M^{\sigma})^{\tau_n}$$

The equation is true because of the definition of $\int_0^{t \wedge \sigma} f dM$ in terms of the stopping times. Therefore, $\int_0^{t \wedge \sigma} f dM = \int_0^t f dM^{\sigma}$ as local martingales. Of course if M is a martingale, this is true also. ■

Bibliography

- [1] **Apostol T. M.**, *Calculus Volume II Second edition*, Wiley 1969.
- [2] **Apostol, T. M.**, *Mathematical Analysis*, Addison Wesley Publishing Co., 1974.
- [3] **Ash, Robert**, *Complex Variables*, Academic Press, 1971.
- [4] **Baker, Roger**, *Linear Algebra*, Rinton Press 2001.
- [5] **Balakrishnan A.V.**, *Applied Functional Analysis*, Springer Verlag 1976.
- [6] **Billingsley P.**, *Probability and Measure*, Wiley, 1995.
- [7] **Buck, R. C.** *Advanced Calculus* 2 edition. McGraw-Hill, 1965.
- [8] **Brézis, H.** *Opérateurs maximaux monotones et semigroupes de contractions dans les espaces de Hilbert*, Math Studies, 5, North Holland, 1973.
- [9] **Brézis, H.**, *Équations et inéquations non linéaires dans les espaces vectoriels en dualité*, Ann. Inst. Fourier (Grenoble) 18 (1968) pp. 115-175.
- [10] **Cheney, E. W.**, *Introduction To Approximation Theory*, McGraw Hill 1966.
- [11] **Chow S.N. and Hale J.K.**, *Methods of bifurcation Theory*, Springer Verlag, New York 1982.
- [12] **Deimling K.** *Nonlinear Functional Analysis*, Springer-Verlag, 1985.
- [13] **Diestal J. and Uhl J.**, *Vector Measures*, American Math. Society, Providence, R.I., 1977.
- [14] **Dontchev A.L.** The Graves theorem Revisited, *Journal of Convex Analysis*, Vol. 3, 1996, No.1, 45-53.
- [15] **Donal O'Regan, Yeol Je Cho, and Yu-Qing Chen**, *Topological Degree Theory and Applications*, Chapman and Hall/CRC 2006.
- [16] **Dunford N. and Schwartz J.T.** *Linear Operators*, Interscience Publishers, a division of John Wiley and Sons, New York, part 1 1958, part 2 1963, part 3 1971.
- [17] **Evans L.C. and Gariepy**, *Measure Theory and Fine Properties of Functions*, CRC Press, 1992.
- [18] **Evans L.C.** *Partial Differential Equations*, Berkeley Mathematics Lecture Notes. 1993.
- [19] **Fitzpatrick P. M.**, *Advanced Calculus a course in Mathematical Analysis*, PWS Publishing Company 1996.
- [20] **Fonesca I. and Gangbo W.** *Degree theory in analysis and applications* Clarendon Press 1995.
- [21] **Gasinski L., Migorski S., and Ochal A.** Existence results for evolutionary inclusions and variational–hemivariational inequalities, *Applicable Analysis* 2014.
- [22] **Gasinski L. and Papageorgiou N.**, *Nonlinear Analysis*, Volume 9, Chapman and Hall, 2006.

- [23] **Gurtin M.** *An introduction to continuum mechanics*, Academic press 1981.
- [24] **Gromes W.** Ein einfacher Beweis des Satzes von Borsuk. *Math. Z.* 178, pp. 399 -400 (1981).
- [25] **Hardy, G.H., Littlewood, J.E. and Polya, G.,** *Inequalities*, Cambridge University Press 1964.
- [26] **Hewitt E. and Stromberg K.** *Real and Abstract Analysis*, Springer-Verlag, New York, 1965.
- [27] **Heinz, E.** An elementary analytic theory of the degree of mapping in n dimensional space. *J. Math. Mech.* 8, 231-247 1959
- [28] **Hobson E.W.,** *The Theory of functions of a Real Variable and the Theory of Fourier's Series V. 1*, Dover 1957.
- [29] **Hocking J. and Young G.,** *Topology*, Addison-Wesley Series in Mathematics, 1961.
- [30] **Horn R. and Johnson C.,** *matrix Analysis*, Cambridge University Press, 1985.
- [31] **Hu S. and Papageorgiou, N.** *Handbook of Multivalued Analysis*, Kluwer Academic Publishers (1997).
- [32] **Karatzas and Shreve,** *Brownian Motion and Stochastic Calculus*, Springer Verlag, 1991.
- [33] **Kato T.** *Perturbation Theory for Linear Operators*, Springer, 1966.
- [34] **Kreyszig E.** *Introductory Functional Analysis With applications*, Wiley 1978.
- [35] **Kuratowski K. and Ryll-Nardzewski C.** A general theorem on selectors, *Bull. Acad. Pol. Sc.*, 13, 397-403.
- [36] **Kuttler K.L.,** *Modern Analysis* CRC Press 1998.
- [37] **Kuttler K. L.,** *Basic Analysis*, Rinton
- [38] **Kuttler K. L.,** *Linear Algebra and Analysis*, web page [Web Page](#)
- [39] **Marsden J. E. and Hoffman J. M.,** *Elementary Classical Analysis*, Freeman, 1993.
- [40] **McShane E. J.** *Integration*, Princeton University Press, Princeton, N.J. 1944.
- [41] **Munkres, James R.,** *Topology A First Course*, Prentice Hall, Englewood Cliffs, New Jersey 1975
- [42] **Natanson I. P.,** *Theory Of Functions Of A Real Variable*, Fredrick Ungar Publishing Co. 1955.
- [43] **Naylor A. and Sell R.,** *Linear Operator Theory in Engineering and Science*, Holt Rinehart and Winston, 1971.
- [44] **Øksendal Bernt** *Stochastic Differential Equations*, Springer 2003.

- [45] **Pettis, B.J.** On integration in vector spaces. Trans. Amer. Math.Soc. 44 277-304 (1938)
- [46] **Prévôt C. and Röckner, A** *Concise Course on Stochastic Partial Differential Equations, Lecture notes in Mathematics, Springer 2007.*
- [47] **Ray W.O.** *Real Analysis*, Prentice-Hall, 1988.
- [48] **Rohatgi V. K.** *An Introduction to Probability Theory and Mathematical Statistics*, John Wiley and Sons, 1976.
- [49] **Rudin W.,** *Principles of Mathematical Analysis*, McGraw Hill, 1976.
- [50] **Rudin W.** *Real and Complex Analysis*, third edition, McGraw-Hill, 1987.
- [51] **Rudin W.** *Functional Analysis*, second edition, McGraw-Hill, 1991.
- [52] **Simon, J** *Compact sets in the space $L^p(0, T; B)$* , Ann. Mat. Pura. Appl. **146**(1987), 65-96.
- [53] **Spanier E.,** *Algebraic Topology*, McGraw Hill 1966.
- [54] **Spivak M.,** *Calculus On Manifolds*, Benjamin 1965.
- [55] **Stromberg, K. R.** *Probability for analysts*, Chapman and Hall, 1994.
- [56] **Stroock D. W.** *Probability Theory An Analytic View*, Second edition, Cambridge University Press, 2011
- [57] **Taylor A. E.** *General Theory of Functions and Integration*, Blaisdell Publishing, 1965
- [58] **Vick, J.,** *homology theory, An Introduction to Algebraic Topology*, Academic Press 1973.
- [59] **Widder, D.** *Advanced Calculus*, second edition, Prentice Hall 1961.
- [60] **Yosida K.** *Functional Analysis*, Springer-Verlag, New York, 1978.

Index

- $(-\infty, \infty]$, 238
- C^k , 193
- C^1 , 192
- C^1 and differentiability, 192
- C_c^∞ , 366
- C_c^m , 366
- F_σ , 253
- G_δ , 534
- G_δ , 253
- L^1
 - approximation, 315
 - weak compactness, 631
- L^p
 - definition, 359
 - density of continuous functions, 363
 - density of simple functions, 362
 - density of smooth functions, 368
 - norm, 359
 - separability, 363
- L^p
 - reflexive, 631
 - weak compactness, 631
- L^1
 - complex vector space, 286
- $L^1(\Omega)$, 285
- L^∞ , 361
- L^p
 - compactness, 371
 - continuity of translation, 365
- $L^p(\Omega)$, 357
- $L^p(\Omega; X)$, 671
- L_{loc}^1 , 366
- $X \times Y$
 - norm, 538
- ε net, 78
- \mathbb{F}^n , 99
- \mathcal{D}^* , 443
- \mathcal{G} , 375
- $\mathcal{L}(X, Y)$, 535
 - Banach space, 535
- π systems, 243
- σ algebra, 237
- a-priori estimates, 167
- absolutely continuous
 - existence of derivative, 346
 - function, 345
 - function and measure, 345
 - integral of derivative, 349
 - integral of the derivative, 346
 - Lipschitz, 347
 - measure, 345
- accumulation point, 71
- adapted, 786, 817
- adjoint linear map, 544
- adjugate, 49
- Alexander subbasis theorem, 507
- algebra, 141
 - Cartesian product, 524
 - measure on algebra, 524
 - recognizing one, 523
- algebra of sets, 523
- approximate identity, 366
- a-priori estimates, 167
- arcwise connected, 90
 - connected, 90, 144
- area measure
 - on manifold, 399
- arithmetic mean, 233
- Arzela Ascoli theorem
 - Banach space, 690
- Ascoli Arzela theorem
 - general form, 87
- at most countable, 61
- axiom of choice, 57, 61
- axiom of extension, 57
- axiom of specification, 57
- axiom of unions, 57
- backwards Holder inequality, 373
- backwards Minkowski inequality, 374
- Baire
 - category, 533, 534
- Baire theorem, 122
- Banach
 - space, 533
- Banach Alaoglu theorem, 556
- Banach space, 106, 205, 359, 533, 711
- Banach Steinhaus theorem, 536
- barycenter, 157
- basis, 101
- basis of a topology, 501
- Bernstein polynomial
 - approximation of derivative, 133
- Besicovitch
 - covering theorem, 119, 317
- Besicovitch covering theorem, 263
- Bessel's inequality, 585, 619

- Binet Cauchy formula, 45
- block matrix, 15
- block multiplication, 15, 17
- Bochner integrable, 655
- Borel
 - measure, 248
- Borel Cantelli lemma, 242, 715
- Borel measure
 - metric space, 248
 - Polish space, 256
- Borel regular, 447
- Borel sets, 252
- Borsuk, 423
- Borsuk Ulam theorem, 425
- Borsuk's theorem, 421
- bounded continuous linear functions, 534
- bounded linear maps, 205
 - continuity, 126
- bounded set, 110
- Brouwer fixed point theorem, 329, 420, 578
 - compact convex set, 163
- Browder's lemma, 178, 278, 619
- Burkholder Davis Gundy
 - inequality, 877
- Burkholder Davis Gundy inequality, 875
- Cantor function, 269
- Cantor set, 268
- Caratheodory extension theorem, 525
- Caratheodory functions, 274
- Caratheodory's procedure, 246
- Cariste fixed point theorem, 172
- Cauchy Schwarz inequality, 104, 575
- Cauchy sequence, 73
- Cayley Hamilton theorem, 53
- central limit theorem, 778
- chain, 68
- chain rule, 188
- change of variables, 467
 - linear map, 330
 - linear maps, 331
- characteristic function, 750
- characteristic polynomial, 53
- Clairaut's theorem, 198
- Clarkson
 - inequalities, 551
- Clarkson inequalities, 552
- Clarkson inequality
 - $p \geq 2$, 548
 - easy one, 549
- closed disk, 110
- closed graph theorem, 539
- closed set, 72, 502
- closed sets
 - limit points, 72
- closure of a set, 74, 75, 503
- coarea formula, 462, 463
- cofactor, 47
- cofactor identity, 201, 414
- column rank, 50
- compact map
 - finite dimensional approximation, 163
- compact set, 76, 504
- compactness
 - closed interval, 109
 - equivalent conditions, 78
- completely separable, 75
- complex
 - measure, 621
- complex valued measurable functions, 284
- components of a vector, 101
- conditional expectation, 781
 - Banach space, 701
 - independence, 784
- connected, 88
 - open balls, 90, 144
- connected component, 89
 - boundary, 410
- connected components, 89
 - equivalence class, 89
 - equivalence relation, 89
 - open sets, 90
- connected set
 - continuous function, 91, 145
 - continuous image, 88
- connected sets
 - intersection, 88
 - intervals, 89
 - real line, 89
 - union, 88
- continuity
 - algebraic properties, 129
 - bounded linear maps, 126
 - coordinate maps, 127
 - uniform, 82
- continuity of translation, 365

- continuity set, 777
- continuous function, 80, 503
 - maximum and minimum, 82
- continuous functions
 - compact support, 290
 - equivalent conditions, 80
- continuous image of compact set, 81
- continuous martingale
 - not of bounded variation, 863
- contraction map, 82
 - fixed point, 208
 - fixed point theorem, 82
- convergence in measure, 310
- convex, 711
 - set, 576
 - sets, 145
- convex
 - functions, 372
- convex combination, 120
- convex function
 - continuous, 741
- convex hull, 120, 155, 711
- convolution, 366, 388
- convolution of measures, 755
- coordinate map, 111
- coordinates, 155
- countable, 61
- covariance matrix, 759
- covariation, 871
- cowlicks, 425
- Cramer's rule, 49
- cylindrical set, 725, 749

- definition of L^p , 359
- definition of a C^k function, 194
- density of \mathcal{G} in L^p , 376
- density of continuous functions in L^p , 363
- derivative
 - chain rule, 188
 - continuity, 193
 - continuity of Gateaux derivative, 193
 - continuous, 187
 - continuous Gateaux derivatives, 191
 - Frechet, 186
 - Gateaux, 189, 191
 - generalized partial, 196
 - higher order, 193
 - matrix, 189
 - partial, 196
 - second, 193
 - well defined, 186
- derivative of inverse, 230
- derivatives, 186
- determinant
 - definition, 41
 - expansion along row, column, 47
 - matrix inverse, 48
 - permutation of rows, 42
 - product, 45
 - row, column operations, 43
 - symmetric definition, 43
 - transpose, 43
- diameter of a set, 96
- differentiable, 186
 - continuous, 187
 - continuous partials, 197
- differentiable map of Lebesgue measurable set, 332
- differential equations
 - dependence on data, 562
 - global existence, 166
- differential forms
 - generalalities, 472
- differentiation
 - Radon measures, 341
- differentiation almost everywhere
 - monotone function, 267
- dimension of a vector space, 102
- Dini derivatives, 265
- directional derivative, 191
- distance, 71
- distribution, 443, 715
- distribution function, 297
- divergence theorem, 402
- dominated convergence
 - generalization, 290
- dominated convergence theorem, 289, 673
- Doob Dynkin lemma, 723
- Doob estimate, 793
- Doob's sub-martingale estimate, 826
- dot product, 103
- dual space, 543
- duality maps, 572
- dyadics, 125
- Dynkin's lemma, 243
 - alternative conditions, 269

- Eberlein Smulian theorem, 560
- Egoroff theorem, 291
- Egoroff theorem, 654
- eigenvalue, 233
- eigenvalues, 53
- Ekeland variational principle, 170
- elementary function, 883, 893
- elementary functions
 - linear space, 893
- elementary matrices, 24
- embedding into its double dual space, 545
- equality of mixed partial derivatives, 200
- equi-integrable, 293
- equicontinuous, 86, 690
- equivalence class, 63
- equivalence relation, 63
- events, 720
- evolution equation
 - continuous semigroup, 610
- exchange theorem, 101
- expectation, 722
 - products of random variables, 723
- exponential growth, 391
- extending off closed set, 711
- extention
 - mapping, 711
- extreme values theorem, 82

- F sigma
 - set, 253
- Faddeyev, 310
- Fatou's lemma, 283
- Fick's law, 406
- filtration, 817
- filtration
 - normal, 817
- finding the inverse, 37
- finite intersection property, 504
- finite measure
 - regularity, 255, 717
- first hitting time, 796
 - closed set, 838
 - open set, 838
- fixed point property, 329, 434
- fixed point theorem
 - Cariste, 172
 - Kakutani, 169
- flip, 224

- Fourier and inverse Fourier transforms, 379
- Fourier series
 - uniform convergence, 571
- Fourier transform
 - L^1 , 381
 - L^2 , 383
 - continuous, 381
 - convolution, 382, 388
 - in \mathcal{G}^* , 378
 - of functions in G , 376
 - convolution, 386
- Fourier transform L^1 , 381
- Fourier transforms
 - polynomial growth, 380
- Frechet derivative, 186
- Fredholm operator
 - Banach space, 564
- Fubini's theorem, 306
 - Bochner integrable functions, 669
 - general product measures, 305
- function, 60
- functions
 - measurable, 237
- fundamental theorem of calculus
 - general Radon measures, 325
 - Radon measures, 324

- G delta, 253
- Gamma function, 372, 452
- gamma function, 204
- Gateaux
 - derivative, 437
- Gateaux derivative, 189, 191
 - continuous, 193
- gauge function, 540
- Gauss Jordan method for inverses, 37
- generalized normal distribution, 762
- geometric mean, 233
- good lambda inequality, 299, 874
- Gram Schmidt process, 109
- Gram Schmidt process., 108
- graph of a linear map, 538
- Green's theorem, 497

- Hahn
 - decomposition, 301
- Hahn Banach theorem, 541
 - complex version, 542

- Hahn decomposition, 645
- Hahn Jordan decomposition, 645
- Hamel basis, 121
- Hardy's inequality, 372
- Hausdorff measures, 445
- Hausdorff
 - maximal principle, 68
- Hausdorff and Lebesgue measure, 452
- Hausdorff dimension, 452
- Hausdorff maximal principle, 507, 540
- Hausdorff measure
 - set of measure 0, 453
- Hausdorff measures, 446
- Hausdorff space, 502
- Hermitian
 - diagonalization, 20
 - non-defective, 20
- Hermitian matrix, 20
- Hessian matrix, 219
- higher order derivative
 - multilinear form, 193
- higher order derivatives, 193
 - implicit function theorem, 212
 - inverse function theorem, 212
- Hilbert Schmidt
 - operator, 592
- Hilbert Schmidt theorem, 586, 665
- Hilbert space, 575
- hitting this before that, 853
- Holder inequality
 - backwards, 546
- Holder space
 - not separable, 567
- Holder spaces, 567
- Holder's inequality, 106, 357
- homeomorphism, 81
- implicit function theorem, 208
 - higher order derivatives, 212
- increasing function
 - existence of the derivative, 437
- independent, 313
- independent events, 720
- independent random vectors, 721
- independent sigma algebras, 721
- indicator function
 - approximation, 315
- inner product space, 575
- inner regular, 253, 717
 - compact sets, 253
- inner regularity, 257
- Integral
 - Riemann and Lebesgue, 531
- integral
 - continuous function, 138
 - decreasing function, 279
 - functions in L^1 , 285
 - linear, 285
- integral over a measurable set, 290
- integrals
 - iterated, 141
- integration
 - with respect to a martingale, 883
- integration by parts, 323
- integration with respect to martingales
 - Ito isometry, 888
- interior point, 71
- intermediate value theorem, 89
- invariance of domain, 230, 424
- inverse, 29, 37
- inverse
 - left right, 38
 - product of matrices, 30
 - row reduced echelon form, 39
- inverse function theorem, 211, 232
 - higher order derivatives, 212
- inverse image, 59
- inverses and determinants, 48
- invertible, 29
- invertible maps, 205
 - different spaces, 206
- isodiametric inequality, 449, 451
- isometric, 682
- iterated integrals, 141
- Ito isometry, 888
- James map, 545
- Jensen's inequality, 741
- Jensens inequality, 784
- Jordan curve theorem, 430
- Jordan separation theorem, 431
- Kakutani fixed point theorem, 169
- Kantorovitch, 635
- Kolmogorov Centsov theorem, 810, 813
- Kolmogorov extension theorem, 527, 719

- Kolmogorov zero one law, 728
- Kolmogorov's inequality, 730
- Kuratowski theorem, 274

- Lagrange multipliers, 216, 217
- Laplace expansion, 46
- Laplace transform, 391, 521
- least squares regression, 203
- Lebesgue
 - decomposition, 303
- Lebesgue integral
 - desires to be linear, 284
 - nonnegative function, 280
 - other definitions, 283
 - simple function, 281
- Lebesgue integral versus Riemann Stieltjes integral, 295
- Lebesgue measurable function
 - approximation with Borel measurable, 315
- Lebesgue measure
 - approximation with Borel sets, 315
 - one dimensional, 258
 - properties, 315
- Lebesgue number, 77, 97
- Lebesgue points, 324
- Lebesgue Stieltjes measure, 250, 257, 345
- left inverse, 38, 39
- Leray Schauder alternative, 165
- Levy theorem, 772
- lim inf, 66
 - properties, 67
- lim sup, 66
 - properties, 67
- lim sup and lim inf, 288
- limit
 - continuity, 185
 - infinite limits, 183
 - point, 71
- limit of a function, 183
- limit of a sequence, 72
 - well defined, 72
- limit point, 183, 502
- limits
 - combinations of functions, 183
 - existence of limits, 66
 - independent random variables, 729
- limits and continuity, 185
- Lindeloff property, 76
- linear
 - not continuous, 535
- linear combination, 30, 44, 100
- linear functional
 - positive, 318
- linear independence, 103
- linear map of measurable set, 330
- linear maps, 15
 - closed, 539
 - continuous, 534
 - equivalent conditions, 534
- linear relationship, 30
- linear relationships
 - row operations, 31
- linear space, 99
- linear transformation
 - defined on a basis, 125
 - dimension of vector space, 125
 - rank m, 220
- linear transformations
 - a vector space, 125
 - sum, 125
- linearly dependent, 100
- linearly independent, 100
- linearly independent set
 - enlarging to a basis, 103
- Lipschitz
 - continuous, 82
 - functions, 437
- Lipschitz function
 - integral of its derivative, 437
- Lipschitz functions, 394
 - of measurable sets, 329
- Lipschitz maps
 - extension, 442
- little o notation, 186
- local martingale, 864
- local maximum, 219
- local minimum, 219
- local submartingale, 864
- localizing sequence, 864
- locally compact, 504
- locally compact , 504
- locally finite, 369, 705
- locally one to one, 233
- lower semicontinuous, 97
- Lusin's theorem, 371

- Lyapunov Schmidt procedure, 563
- manifold
 - orientable, 396
- map
 - C^1 , 225
 - primitive and flips, 225
- martingale, 741
 - quadratic variation, 866
- martingales
 - equiintegrable, 851
- matrix
 - left inverse, 49
 - lower triangular, 49
 - right inverse, 49
 - right, left inverse, 49
 - row, column, determinant rank, 50
 - upper triangular, 49
- matrix
 - block multiplication, 17
 - inverse, 29, 37
 - more columns than rows, 35
 - non zero kernel, 35
 - partitioned, 17
 - Schur's theorem, 19
- matrix multiplication
 - block, 17
- maximal chain, 68
- maximal estimate
 - real sub-martingales, 801
- maximal function
 - Radon measures, 324
- McShane's lemma, 637
- mean value inequality, 208
- mean value theorem, 208
 - Cauchy, 137
- measurability
 - limit of simple functions, 239
- measurable, 245
 - complex valued, 284
 - equivalent formulations, 238
 - linear combinations, 284
 - multifunction, 271
 - multifunctions, 271
- measurable complex functions
 - simple functions, 288
- measurable functions, 237
 - approximation, 240
 - pointwise limit, 237
 - pointwise limits, 651
 - simple functions, 238
- measurable into $(-\infty, \infty]$, 238
- measurable representative, 678
- measurable selection, 272
- measurable sets, 245
- measure, 241
 - inner regular, 253
 - on an algebra, 524
 - outer regular, 253
 - properties, 242
 - vector, 621
- measure space
 - completion, 260
 - regular, 362
- measures
 - complex, 621
 - decreasing sequences of sets, 242
 - increasing sequences of sets, 242
 - regularity, 253, 256
 - tight, 766
 - weak convergence, 770
- measures from outer measures, 246
- metric, 71
 - properties, 71
- metric space, 71
 - compact sets, 78
 - complete, 74
 - completely separable, 75
 - open set, 71
 - separable, 75
- min max theorem, 147
- min-max theorem, 147
- Minkowski functional, 571
- Minkowski inequality, 358
 - integrals, 360
- Minkowski inequality
 - backwards, 546
- Minkowski theorem
 - for integrals, 373
- Minkowski's inequality, 360
- minor, 47
- mixed partial derivatives, 198
- modification, 807
- mollifier, 366
- monotone convergence theorem, 282
- monotone function

- differentiable, 267
- multi - index, 130
- multi-index, 194, 375
- multi-index notation, 194
- multifunction
 - strongly measurable, 273
- multifunctions, 271
 - measurability, 271
- Muntz theorem, 153
- negative part, 285
- Neuman series, 205
- Newton's method, 201
- no retract onto boundary of ball, 421
- non equal mixed partials
 - example, 199
- norm
 - p norm, 106
- normal, 405, 759, 773
- normal filtration, 817
- normal topological space, 503
- Normed linear space, 106
- nowhere differentiable functions, 569
- nuclear operator, 590
- one dimensional Stieltjes measure, 294
- one point compactification, 506
- open ball, 71
 - open set, 71
- open cover, 76, 504
- open mapping theorem, 536
- open set, 71
- open sets, 501
 - countable basis, 75
- operator norm, 534
- optional sampling theorem, 798
- ordered
 - partial, 68
 - totally ordered, 68
- oriented atlas, 396
- orthonormal, 108
- orthonormal set, 584
- outer measure, 244
 - measurable, 245
 - total variation, 622
- outer measure on real line, 250
- outer regular, 253, 717
 - G delta and F sigma sets, 255
- outer regularity, 257
- p norms, 107
- paracompact space, 705
- parallelogram identity, 619
- partial derivatives, 189, 196
 - continuous, 197
- partial order, 68
- partially ordered set, 68
- partition
 - one dimension, 348
- partition of unity, 92, 512, 709
 - metric space, 708, 710
- partitioned matrix, 17
- permutation, 41
- permutation matrix, 24
- pi systems, 243
- pivot column, 33
- Plancherel theorem, 383, 384
- point of density, 326
- pointwise compact, 86
- pointwise convergence, 85
- polar coordinates
 - integral, 340
- polar decomposition, 625
- Polish space, 76, 252
- polynomial, 375
 - in many variables, 130
- polynomials, 130
- positive, 595
- positive linear functional, 315, 516
 - measure, 318
- positive part, 285
- positive self adjoint
 - products, 595
 - roots, 596
- power set, 57
- precompact, 504, 521
- predictable, 823
- preserving distance, 22
- primitive, 224
- prior, 796
- prior sets, 836
- probability distribution function, 250
- probability space, 715
- product measure, 306
 - regular, 306
- product of matrices

- inverse, 30
- product space
 - norm, 538
- product topology, 504
- products of Borel sets, 259, 322
- progressively measurable, 817
- progressively measurable
 - integral of, 819
- progressively measurable version, 821
- projection in Hilbert space, 577
- Prokhorov's theorem, 770
- quadratic variation
 - convergence in probability, 904
 - fantastic properties, 887
- Rademacher's theorem, 441
- Radon measure, 253, 263, 362, 363
- Radon Nikodym
 - Radon measures, 343
 - theorem, 302
- Radon Nikodym derivative, 302
- Radon Nikodym property, 680
- Radon Nikodym theorem
 - Radon Measures, 343
- random variable, 715
 - distribution measure, 256
- random vector, 715
 - independent, 736
- rank, 35
- rank of a matrix, 50
- rank theorem, 221
- rational function, 130
- real and imaginary parts, 284
- recognizing a martingale
 - stopping times, 859
- refinement of a cover, 705
- reflexive Banach Space, 546
- reflexive Banach space, 631
 - weak compactness, 560
- regular, 253
- regular measure space, 362
- regular topological space, 502
- regular values, 410
- resolvent, 605
- retract, 329
 - Banach space, 712
 - closed and convex set, 712
 - fixed point property, 329
- reverse sub-martingale, 802
- Reynolds
 - transport formula, 494
- Riemann integral, Lebesgue integral, 294
- Riesz F., 374
- Riesz map, 579
- Riesz representation theorem, 684
 - Hilbert space, 579
 - locally compact Hausdorff space, 516
- Riesz Representation theorem
 - $C(X)$, 642
- Riesz representation theorem L^p
 - σ finite case, 630, 687
 - finite measures, 626
- Riesz representation theorem for L^1
 - finite measures, 629
- right inverse, 37–39
- right polar factorization, 21
- row equivalent, 33
- row operations, 24
- row operations
 - linear relationships, 31
- row rank, 50
- row reduced echelon form, 32
- row reduced echelon form
 - existence, 33
 - uniqueness, 34
- Russell's paradox, 59
- saddle point, 147
- Sard's theorem, 332, 412
- scalars, 99
- Schaefer fixed point theorem, 165
- Schauder fixed point
 - approximate fixed point, 164
- Schauder fixed point theorem, 165, 166
- Schroder Bernstein theorem, 60
- Schur's theorem, 19
- Schwartz class, 386
- second derivative, 193
- second derivative test, 219, 220
- sections of open sets, 196
- self adjoint, 595
- semigroup, 600
 - adjoint, 614
 - contraction
 - bounded, 603

- generator, 602
- growth estimate, 603
- Hille Yosida theorem, 606
- strongly continuous, 602
- separability of $C(H)$, 769
- separable metric space
 - Lindeloff property, 76
- separated sets, 88
- separation theorem, 572
- sequence, 72
 - Cauchy, 73
 - subsequence, 73
- sequential weak* compactness, 558
- sequentially compact set, 77
- series
 - double sum, 64
- set
 - F sigma, 253
- set valued functions
 - measurability, 271
- sets, 57
 - G delta, 253
- sgn, 39
 - uniqueness, 41
- sigma algebra, 237
- sign of a permutation, 41
- signed measure, 300
 - Hahn decomposition, 301
- signed measures
 - Hahn decomposition, 300
- simple functions, 647
 - approximation, 238
- singular values, 410
- slicing measures, 308
- smooth manifold, 396
- Sobolev Space
 - embedding theorem, 391
 - equivalent norms, 391
- Sobolev spaces, 391
- space of continuous martingales, 856
 - Hilbert space, 879
- span, 44, 100
- Sperner's lemma, 160
- Steiner symetrization, 450
- Stirling's formula, 353
- stochastic integral
 - as Stieltjes integral, 891
 - elementary function, 893
 - linear, 893
- stochastic process, 786, 807
 - descriptions, 807
- Stone Weierstrass theorem, 510
- Stone's theorem, 707
- strong law of large numbers, 792
- stopped martingale, 862, 864
- stopped process, 853
- stopped sub-martingale, 845
- stopped submartingale, 864
- stopping time, 795, 831, 835
- stopping times
 - conditional expectation, 797
- strict convexity, 572
- strictly convex
 - norm, 147
- strong law of large numbers, 733, 805
- strongly measurable, 647
 - inverse images open sets, 648
- sub-martingale, 741, 786
- sub-martingale convergence theorem, 789
 - continuous case, 850
- subbasis, 508
- subspace, 100
- sums
 - independent random variables, 729
 - independent variables, 731, 790
- supermartingale, 786
- support of a function, 92, 290, 511
- symmetric derivative
 - existence, 342
 - measurable, 342
 - Radon measure, 341
 - upper and lower, 341
- symmetric domain
 - degree, 423
- symmetric matrix, 20
- tail event, 728
- Taylor formula, 217
- Taylor's formula, 218
- Taylor's theorem, 218
- Tietze extension theorem, 136
- tight, 766
- topological space, 501
- total variation, 348, 622
- totally bounded, 78, 86
- totally ordered, 68

- trace, 591
- trajectories, 807
- translation invariant, 258
- triangle inequality, 106, 358
- triangulated, 155
- triangulation, 155
- trivial, 100

- uniform boundedness, 536
- uniform boundedness theorem, 536
- uniform continuity, 82
- uniform continuity and compactness, 82
- uniform convergence, 85
- uniform convergence and continuity, 85
- uniform convexity, 546, 572
- uniformly bounded, 690
- uniformly convex, 546
- uniformly integrable, 292, 372
- uniqueness of limits, 183
- upcrossing, 741, 787
- upper semicontinuity
 - set valued map, 167
- upper semicontinuous, 97
- upper semicontinuous composition, 168
- Urysohn's lemma, 505, 506

- variance, 759
- variational inequality, 577
- vector
 - measure, 621
- vector measures, 679
- vector space
 - dimension, 102
- vector space axioms, 99
- vector valued function
 - limit theorems, 183
- vectors, 99
- version, 807
- Vitali
 - convergence theorem, 292
- Vitali
 - convergence theorem, 673
- Vitali convergence theorem, 293
- Vitali cover, 263
- Vitali covering theorem, 114
- volume of unit ball, 452

- weak and weak* topologies, 554
- weak convergence, 573
 - measures, 770
- weak convergence of measures, 772
- weak derivative, 444
- weak topology, 555
- weak* measurable, 653
- weak* topology, 555
 - metric space, 557
- weakly measurable, 647
- wedge product, 475
 - algebraic properties, 476
- Weierstrass
 - Stone Weierstrass theorem, 143
- Weierstrass approximation
 - estimate, 130
- well ordered sets, 70

- Young's inequality, 644