
Topics in Undergraduate Mathematics

Ryan Joo

Email: ryanjooruian18@gmail.com

Copyright © 2025 by Ryan Joo.

This is (still!) an incomplete draft. Please send corrections and comments to ryanjooruian18@gmail.com, or pull-request at <https://github.com/Ryanjoo18/undergrad-maths>.

Last updated March 19, 2025.

Preface

I began writing this book during the mid-year break in 2023, which was when I had a bit of free time to read up on undergraduate mathematics, and also started learning \LaTeX to write this book.

The objective of this book is to serve as a compilation of essential topics at the undergraduate level (as well as serving as my personal notes). The book covers (or *aims* to cover) the following topics:

1. Abstract algebra, which follows [DF04].
2. Linear algebra, which follows [Ax124].
3. Real analysis, which follows [Rud76; Apo57]; multivariable analysis, which follows [Spi65; Rud76].
4. General topology, which follows [Mun18].
5. Measure theory and complex analysis, which follows [Rud87].
6. Functional analysis, which follows [Rud91].
7. Stochastic analysis
8. Differential geometry

Prerequisites

This book is written such that it is accessible to high school students. No formal prerequisites are required, although some experience with proofs may be helpful.

Presentation

This book follows the typical style of “Definition”, “Theorem”, etc.

For ease of reference, important terms are *coloured* when first defined, and are included in the glossary; less important terms are *italicised* when first defined, and are not included in the glossary.

Note on Problem Solving

Mathematics is about problem solving. In [Pól45], George Pólya outlined the following problem solving cycle.

1. **Understand the problem**

Ask yourself the following questions:

- Do you understand all the words used in stating the problem?
- Is it possible to satisfy the condition? Is the condition sufficient to determine the unknown? Or is it insufficient? Or redundant? Or contradictory?
- What are you asked to find or show? Can you restate the problem in your own words?
- Draw a figure. Introduce suitable notation.
- Is there enough information to enable you to find a solution?

2. Devise a plan

A partial list of heuristics – good rules of thumb to solve problems – is included:

- Guess and check
- Look for a pattern
- Make an orderly list
- Draw a picture
- Eliminate possibilities
- Solve a simpler problem
- Use symmetry
- Use a model
- Consider special cases
- Work backwards
- Use direct reasoning
- Use a formula
- Solve an equation
- Be ingenious

3. Execute the plan

This step is usually easier than devising the plan. In general, all you need is care and patience, given that you have the necessary skills. Persist with the plan that you have chosen. If it continues not to work discard it and choose another. Don't be misled, this is how mathematics is done, even by professionals.

- Carrying out your plan of the solution, check each step. Can you see clearly that the step is correct? Can you prove that it is correct?

4. Check and expand

Pólya mentions that much can be gained by taking the time to reflect and look back at what you have done, what worked, and what didn't. Doing this will enable you to predict what strategy to use to solve future problems.

Look back reviewing and checking your results. Ask yourself the following questions:

- Can you check the result? Can you check the argument?
- Can you derive the solution differently? Can you see it at a glance?
- Can you use the result, or the method, for some other problem?

Building on Pólya's problem solving strategy, Schoenfeld [Sch92] came up with the following framework for problem solving, consisting of four components:

1. **Cognitive resources:** the body of facts and procedures at one's disposal.
2. **Heuristics:** 'rules of thumb' for making progress in difficult situations.
3. **Control:** having to do with the efficiency with which individuals utilise the knowledge at their disposal. Sometimes, this is referred to as metacognition, which can be roughly translated as 'thinking about one's own thinking'.
 - (a) These are questions to ask oneself to monitor one's thinking.
 - What (exactly) am I doing? [Describe it precisely.] Be clear what I am doing NOW. Why am I doing it? [Tell how it fits into the solution.]
 - Be clear what I am doing in the context of the BIG picture – the solution. Be clear what I am going to do NEXT.

(b) Stop and reassess your options when you

- cannot answer the questions satisfactorily [probably you are on the wrong track]; OR
- are stuck in what you are doing [the track may not be right or it is right but it is at that moment too difficult for you].

(c) Decide if you want to

- carry on with the plan,
- abandon the plan, OR
- put on hold and try another plan.

4. **Belief system:** one's perspectives regarding the nature of a discipline and how one goes about working on it.

Contents

I Preliminary Topics	1
1 Mathematical Reasoning and Logic	2
1.1 Zeroth-order Logic	2
If, only if	4
If and only if, iff	5
1.2 First-order Logic	6
1.3 Methods of Proof	7
Proof by Contradiction	8
Proof of Existence and Uniqueness	9
Proof by Mathematical Induction	12
Pigeonhole Principle	16
Exercises	17
2 Set Theory	23
2.1 Basics of Naive Set Theory	23
Definitions and Notations	23
Algebra of Sets	25
2.2 Relations	28
Definition and Examples	28
Properties of Relations	28
Equivalence Relations	29
Axiom of Choice and Its Equivalences	32
2.3 Functions	33
Definitions	33
Injectivity, Surjectivity, Bijectivity	33
Images and Pre-images	34
Composition	35
Invertibility	38
2.4 Cardinality	41
Exercises	45

II	Abstract Algebra	48
3	Groups	49
3.1	Definition and Properties	49
	Examples	51
	Subgroups	53
	Cyclic Groups	54
	Order	56
3.2	Cosets	57
	Lagrange's Theorem	58
	Counting Principle	59
	Normal Subgroups, Quotient Groups	61
3.3	Homomorphisms and Isomorphisms	63
	Kernel and Image	64
	Isomorphism Theorems	65
3.4	Group Actions	67
	Conjugation	69
	Sylow's Theorem	69
3.5	Group Product, Finite Abelian Groups	70
	Exercises	71
III	Linear Algebra	72
4	Finite Dimensional Vector Spaces	73
4.1	Definition of Vector Space	73
4.2	Subspaces	76
4.3	Span and Linear Independence	80
4.4	Bases	84
4.5	Dimension	87
	Exercises	89
5	Linear Maps	93
5.1	Vector Space of Linear Maps	93
5.2	Kernel and Image	96
	Fundamental Theorem of Linear Maps	98
5.3	Matrices	101
	Representing a Linear Map by a Matrix	101

Addition and Scalar Multiplication of Matrices	102
Matrix Multiplication	104
Rank of a Matrix	107
5.4 Invertibility and Isomorphism	109
Invertibility	109
Isomorphism	111
Linear Maps Thought of as Matrix Multiplication	113
Change of Basis	115
5.5 Products and Quotients of Vector Spaces	118
Products of Vector Spaces	118
Quotient Spaces	121
5.6 Duality	125
Dual Space and Dual Map	125
Kernel and Image of Dual of Linear Map	128
Matrix of Dual of Linear Map	130
Exercises	131
6 Polynomials	134
6.1 Definitions	134
6.2 Zeros of Polynomials	135
6.3 Division Algorithm for Polynomials	137
6.4 Factorisation of Polynomials over \mathbb{C}	138
6.5 Factorisation of Polynomials over \mathbb{R}	139
7 Eigenvalues and Eigenvectors	141
7.1 Invariant Subspaces	141
Eigenvalues	141
Polynomials Applied to Operators	143
7.2 The Minimal Polynomial	145
Existence of Eigenvalues on Complex Vector Spaces	145
Eigenvalues and the Minimal Polynomial	147
Eigenvalues on Odd-Dimensional Real Vector Spaces	150
7.3 Upper-Triangular Matrices	152
7.4 Diagonalisable Operators	155
Diagonal Matrices	155
Conditions for Diagonalisability	156
Gershgorin Disk Theorem	157

7.5	Commuting Operators	158
	Exercises	159
8	Inner Product Spaces	160
8.1	Inner Products and Norms	160
	Inner Products	160
	Norms	162
8.2	Orthonormal Bases	165
	Orthonormal Bases	165
	Gram–Schmidt Procedure	168
	Linear Functionals on Inner Product Spaces	170
8.3	Orthogonal Complements and Minimisation Problems	171
	Orthogonal Complements	171
	Minimisation Problems	173
	Pseudoinverse	174
	Exercises	175
9	Operators on Inner Product Spaces	176
9.1	Self-Adjoint and Normal Operators	176
	Adjoints	176
	Self-Adjoint Operators	177
	Normal Operators	178
9.2	Spectral Theorem	179
	Real Spectral Theorem	179
	Complex Spectral Theorem	179
9.3	Positive Operators	180
9.4	Isometries, Unitary Operators, and Matrix Factorisation	181
	Isometries	181
	Unitary Operators	181
	QR Factorisation	181
	Cholesky Factorisation	181
9.5	Singular Value Decomposition	182
	Singular Values	182
	SVD for Linear Maps and for Matrices	182
9.6	Consequences of Singular Value Decomposition	183
	Norms of Linear Maps	183
	Approximation by Linear Maps with Lower-Dimensional Range	183

Polar Decomposition	183
Operators Applied to Ellipsoids and Parallelepipeds	183
Volume via Singular Values	183
Properties of an Operator as Determined by Its Eigenvalues	183
IV Real Analysis	184
10 Real and Complex Number Systems	185
10.1 Ordered Sets and Boundedness	185
Definitions	185
Least-upper-bound Property	187
Properties of Suprema and Infima	188
Ordered Fields	190
10.2 Real Numbers	191
Problems with \mathbb{Q}	191
Real Field	192
Properties of \mathbb{R}	198
Extended Real Number System	202
10.3 Complex Field	203
10.4 Euclidean Space	207
Exercises	209
11 Basic Topology	211
11.1 Metric Spaces	212
Definitions and Examples	212
Balls and Boundedness	214
Open and Closed Sets	216
Interior, Closure, Boundary	219
Limit Points	221
11.2 Compactness	224
Definitions and Properties	224
Heine–Borel Theorem	228
Bolzano–Weierstrass Theorem	231
Cantor’s Intersection Theorem	232
Sequential Compactness	234
11.3 Perfect Sets	235
Definition and Uncountability	235

Cantor Set	236
11.4 Connectedness	238
Path Connectedness	240
11.5 Separable Spaces	241
11.6 Baire Category Theorem	242
Exercises	244
12 Numerical Sequences and Series	246
12.1 Sequences	246
Convergence	246
Subsequences	253
Cauchy Sequences	255
Monotonic Sequences	258
Limit Superior and Inferior	259
12.2 Series	262
Convergence Tests	263
Summation by Parts	270
Addition and Multiplication of Series	272
Rearrangements	274
Exercises	276
13 Continuity	280
13.1 Limit of Functions	280
13.2 Continuous Functions	282
Continuity and Pre-images of Open or Closed Sets	285
Continuity and Compactness	286
Bolzano's Theorem	288
Continuity and Connectedness	289
13.3 Uniform Continuity	290
13.4 Discontinuities	293
13.5 Monotonic Functions	295
13.6 Lipschitz Continuity	297
13.7 Infinite Limits and Limits at Infinity	299
Exercises	300
14 Differentiation	301
14.1 The Derivative of A Real Function	301

Definitions and Properties	301
Derivatives of Higher Order	305
14.2 Mean Value Theorems	306
14.3 Continuity of Derivatives	308
14.4 L'Hopital's Rule	309
14.5 Taylor's Theorem	311
Exercises	313
15 Riemann–Stieltjes Integral	314
15.1 Definition of Riemann–Stieltjes Integral	314
Notation and Preliminaries	314
Defining the Integral	317
Useful Identities	319
15.2 Properties of the Integral	323
15.3 Integration and Differentiation	330
15.4 Integration of Vector-valued Functions	332
15.5 Rectifiable Curves	334
Exercises	336
16 Sequences and Series of Functions	337
16.1 Pointwise Convergence	337
16.2 Uniform Convergence	339
16.3 Properties of Uniform Convergence	342
Uniform Convergence and Continuity	342
Uniform Convergence and Integration	345
Uniform Convergence and Differentiation	347
16.4 Equicontinuous Families of Functions	350
16.5 Stone–Weierstrass Approximation Theorem	354
Weierstrass's Version	354
Algebra of Functions	356
The Theorem	356
Exercises	357
17 Some Special Functions	358
17.1 Power Series	358
Exponential and Logarithmic Functions	367
Trigonometric Functions	369

17.2 Algebraic Completeness of the Complex Field	370
17.3 Fourier Series	371
17.4 Gamma Function	378
Exercises	385
V Multivariable Analysis	386
18 Functions of Several Variables	387
18.1 Linear Transformations	387
18.2 Differentiation	387
The Derivative	387
Partial Derivatives	391
Gradients, Curves, and Directional Derivatives	393
The Jacobian	395
18.3 Continuity and The Derivative	396
18.4 Inverse and Implicit Function Theorems	397
Inverse Function Theorem	397
Implicit Function Theorem	398
18.5 Derivatives of Higher Order	399
18.6 Differentiation of Integrals	399
Exercises	400

I

Preliminary Topics

1 Mathematical Reasoning and Logic

Summary

- Basic logic.
- Common methods of proof.

§1.1 Zeroth-order Logic

A **proposition** is a sentence which has exactly one truth value, i.e. it is either true or false, but not both and not neither. A proposition is denoted by uppercase letters such as P and Q . If the proposition P depends on a variable x , it is sometimes helpful to denote it by $P(x)$.

We can do some algebra on propositions:

- equivalence**, denoted by $P \equiv Q$, means P and Q are logically equivalent statements;
- conjunction**, denoted by $P \wedge Q$, means “ P and Q ”;
- disjunction**, denoted by $P \vee Q$, means “ P or Q ”;
- negation**, denoted by $\neg P$, means “not P ”.

Here are some useful properties when handling logical statements. You can easily prove all of them using truth tables.

Lemma 1.1.

(i) *Double negation law:*

$$P \equiv \neg(\neg P)$$

(ii) *Commutative laws:*

$$P \wedge Q \equiv Q \wedge P$$

$$P \vee Q \equiv Q \vee P$$

(iii) *Associative laws:*

$$(P \wedge Q) \wedge R \equiv P \wedge (Q \wedge R)$$

$$(P \vee Q) \vee R \equiv P \vee (Q \vee R)$$

(iv) *Idempotent laws:*

$$P \wedge P \equiv P$$

$$P \vee P \equiv P$$

(v) *Distributive laws:*

$$P \wedge (Q \vee R) \equiv (P \wedge Q) \vee (P \wedge R)$$

$$P \vee (Q \wedge R) \equiv (P \vee Q) \wedge (P \vee R)$$

(vi) *Absorption laws:*

$$P \vee (P \wedge Q) \equiv P$$

$$P \wedge (P \vee Q) \equiv P$$

(vii) *de Morgan's laws:*

$$\neg(P \vee Q) \equiv (\neg P \wedge \neg Q)$$

$$\neg(P \wedge Q) \equiv (\neg P \vee \neg Q)$$

Remark. Notice that because of the associative laws we can leave out parentheses in statements of the forms $P \wedge Q \wedge R$ and $P \vee Q \vee R$ without ambiguity, because the two possible ways of filling in the parentheses are equivalent.

Statements that are always true are called *tautologies*; for instance $P \vee \neg P$. Similarly, statements that are always false are called *contradictions*; for instance $P \wedge \neg P$.

We can now state a few more useful laws involving tautologies and contradictions.

Lemma 1.2.

(i) *Tautology laws: if Q is a tautology, then*

$$P \wedge Q \equiv P$$

$$P \vee Q \text{ is a tautology}$$

$$\neg Q \text{ is a contradiction}$$

(ii) *Contradiction laws: if Q is a contradiction, then*

$$P \vee Q \equiv P$$

$$P \wedge Q \text{ is a contradiction}$$

$$\neg Q \text{ is a tautology}$$

If, only if

Implication is denoted by $P \implies Q$, which means “ P implies Q ”, i.e. if P holds then Q also holds. It is equivalent to saying “If P then Q ”. $P \implies Q$ is known as a *conditional statement*, where P is known as the *hypothesis* (or *premise*) and Q is known as the *conclusion*.

The only case when $P \implies Q$ is false is when the hypothesis P is true and the conclusion Q is false.

Statements of this form are probably the most common, although they may sometimes appear quite differently.

The following all mean the same thing:

- (i) if P then Q ;
- (ii) P implies Q ;
- (iii) P only if Q ;
- (iv) P is a sufficient condition for Q ;
- (v) Q is a necessary condition for P .

Given $P \implies Q$,

- its **converse** is $Q \implies P$; both are not logically equivalent;
- its **inverse** is $\neg P \implies \neg Q$, i.e. the hypothesis and conclusion of the statement are both negated; both are not logically equivalent;
- the **contrapositive** is $\neg Q \implies \neg P$; both are logically equivalent.

To prove $P \implies Q$,

1. assume that P holds,
2. deduce, through some logical steps, that Q holds.

Alternatively, we can prove the contrapositive: assume that Q does not hold, then show that P does not hold.

If and only if, iff

Bidirectional implication is denoted by $P \iff Q$, which means both $P \implies Q$ and $Q \implies P$; $P \iff Q$ is known as a *biconditional statement*. We can read this as “ P if and only if Q ”. The letters “iff” are also commonly used to stand for “if and only if”.

$P \iff Q$ is true exactly when P and Q have the same truth value.

These statements are usually best thought of separately as “if” and “only if” statements. To prove $P \iff Q$, prove the statement in both directions:

1. prove $P \implies Q$, and
2. prove $Q \implies P$.

Remember to make very clear, both to yourself and in your written proof, which direction you are doing.

§1.2 First-order Logic

The **universal quantifier** is denoted by \forall , which means “for all” or “for every”. A *universal statement* takes the form $\forall x \in X, P(x)$.

The **existential quantifier** is denoted by \exists , which means “there exists”. An *existential statement* takes the form $\exists x \in X, P(x)$, where X is known as the *domain*.

Lemma 1.3 (de Morgan’s laws).

$$\neg [\forall x \in X, P(x)] \equiv \exists x \in X, \neg P(x)$$

$$\neg [\exists x \in X, P(x)] \equiv \forall x \in X, \neg P(x)$$

To prove a statement of the form $\forall x \in X, P(x)$,

1. Start with “*let $x \in X$ be given*” to address the quantifier with an arbitrary x (this will prove the statement for all $x \in X$).
2. Show that $P(x)$ is true.

Consider statements of the form $\forall x \in X, P(x) \implies Q(x)$; we say that the statement is *vacuously true* if $P(x)$ is false for all $x \in X$.

To prove a statement of the form $\exists x \in X, P(x)$, there is not such a clear steer about how to continue:

- you can construct such an x with the desired properties (constructive proof);
- you can demonstrate logically that such an x must exist because of some earlier assumption (non-constructive proof);
- you can suppose that such an x does not exist, and consequently arrive at some inconsistency (proof by contradiction).

Remark. Read from left to right, and as new elements or statements are introduced they are allowed to depend on previously introduced elements but cannot depend on things that are yet to be mentioned.

Remark. To avoid confusion, it is a good idea to keep to the convention that the quantifiers come first, before any statement to which they relate.

§1.3 *Methods of Proof*

What is a *proof*? Informally, we will define a mathematical proof to be a logical argument that establishes the truth of a mathematical statement. A typical proof proceeds as follows:

1. Start with the given hypotheses.
2. Apply rules of inferences (logical deduction) to get new statements.
3. Repeat Step 2 until we reach the desired conclusion.

We first present some straightforward methods of proof:

- A **direct proof** of $P \implies Q$ is a series of valid arguments that start with the hypothesis P and end with the conclusion Q .

$$P \implies \dots \implies Q$$

- A **proof by contrapositive** of $P \implies Q$ is to prove instead $\neg Q \implies \neg P$.
- A **disproof by counterexample** is to provide a counterexample to disprove a statement, which makes the negation of the statement true.

Thus, to disprove $P \implies Q$, the counterexample makes the hypothesis P true, and the conclusion Q false. Likewise, to disprove $\forall x \in X, P(x)$, we prove its negation $\exists x \in X, \neg P(x)$, i.e., find $a \in X$ such that $P(a)$ is false.

In seeking counterexamples, it is a good idea to keep the cases you consider simple, rather than searching randomly. It is often helpful to consider “extreme” cases; for example, something is zero, a set is empty, or a function is constant.

- A **proof by cases** is to first dividing the situation into cases which exhaust all the possibilities, and then show that the statement follows in all cases.

Proof by Contradiction

To *prove by contradiction*,

1. Assume P is false, i.e., $\neg P$ is true (to prove $P \implies Q$ by contradiction, suppose $P \wedge \neg Q$).
2. Show, through some logical reasoning, that this leads to a contradiction or inconsistency.

We may arrive at something that contradicts the hypothesis P , or something that contradicts the initial supposition that Q is not true, or we may arrive at something that we know to be universally false.

We illustrate this method of proof using a classic example.

Example 1.4 (Irrationality of $\sqrt{2}$). Prove that $\sqrt{2}$ is irrational.

Proof. We prove by contradiction. Suppose otherwise, that $\sqrt{2}$ is rational. Then $\sqrt{2} = \frac{a}{b}$ for some $a, b \in \mathbb{Z}$, $b \neq 0$, a, b coprime.

Squaring both sides gives

$$a^2 = 2b^2.$$

Since RHS is even, LHS must also be even. Hence it follows that a is even. Let $a = 2k$ where $k \in \mathbb{Z}$. Substituting $a = 2k$ into the above equation and simplifying it gives us

$$b^2 = 2k^2.$$

This means that b^2 is even, from which follows again that b is even. This contradicts the assumption that a and b coprime, so we are done. \square

Example 1.5 (Euclid). Prove that there are infinitely many prime numbers.

Proof. Suppose otherwise, that only finitely many prime numbers exist. List them as p_1, \dots, p_n . Consider the number

$$N = p_1 p_2 \cdots p_n + 1.$$

Note that N is divisible by a prime p , yet is coprime to p_1, \dots, p_n . Therefore, p does not belong to our list of all prime numbers, a contradiction. \square

Proof of Existence and Uniqueness

To prove existential statements, we can adopt two approaches:

1. *Constructive proof* (direct proof)

To prove statements of the form $\exists x \in X, P(x)$, find or construct *a specific example* for x . To prove statements of the form $\forall y \in Y, \exists x \in X, P(x, y)$, construct example for x in terms of y (since x is dependent on y).

In both cases, you have to justify that your example x

- (a) belongs to the domain X , and
- (b) satisfies the condition P .

2. *Non-constructive proof* (indirect proof)

Use when specific examples are not easy or not possible to find or construct. Make arguments why such objects have to exist. May need to use proof by contradiction. Use definition, axioms or results that involve existential statements.

To *prove uniqueness* (after proving existence), we can either

- assume $\exists x, y \in X$ such that $P(x) \wedge P(y)$ is true, then show $x = y$, or
- assume that $\exists x, y \in X$ are distinct such that $P(x) \wedge P(y)$, then derive a contradiction.

We sometimes use $\exists!$ to mean “there exists a unique”.

Example 1.6. Prove that we can find 100 consecutive positive integers which are all composite numbers.

Proof. We proceed by constructive proof; we will construct integers $n, n + 1, n + 2, \dots, n + 99$, all of which are composite.

Claim. $n = 101! + 2$.

Then n has a factor of 2 and hence is composite. Similarly, $n + k = 101! + (k + 2)$ has a factor $k + 2$ and hence is composite for $k = 1, 2, \dots, 99$.

Hence the existential statement is proven. □

Example 1.7. Prove that for all $p, q \in \mathbb{Q}$ with $p < q$, there exists $x \in \mathbb{Q}$ such that $p < x < q$.

Proof. We prove by construction; we want to construct x in terms of p and q , which fulfils the required condition.

Claim. $x = \frac{p + q}{2}$.

Evidently $x \in \mathbb{Q}$. Since $p < q$,

$$x = \frac{p + q}{2} < \frac{q + q}{2} = q \implies x < q.$$

Similarly,

$$x = \frac{p + q}{2} > \frac{p + p}{2} = p \implies p < x.$$

Remark. There are two parts to prove: 1) x satisfies the given statement 2) x is within the domain (for this question we do not have to prove x is rational since \mathbb{Q} is closed under addition).

□

Example 1.8. Prove that for all rational numbers p and q with $p < q$, there is an irrational number r such that $p < r < q$.

Proof. We prove this by construction. Similarly, our goal is to find an irrational r in terms of p and q .

Note that we cannot simply take $r = \frac{p+q}{2}$; a simple counterexample is the case $p = -1, q = 1$ where $r = 0$ is clearly not irrational.

Since p lies in between p and q , let $r = p + c$ where $0 < c < q - p$. Since $c < q - p$, we have $c = \frac{q-p}{k}$ for some $k > 1$; to make c irrational, we take k to be irrational.

Claim. $r = p + \frac{q-p}{\sqrt{2}}$.

We shall show that (i) $p < r < q$, and (ii) r is irrational.

(i) Since $q - p > 0$, $\frac{q-p}{\sqrt{2}} > 0$ so $r = p + \frac{q-p}{\sqrt{2}} > p + 0 = p$.

$\frac{q-p}{\sqrt{2}} < q - p$ so $r < p + (q - p) = q$.

(ii) We prove by contradiction. Suppose r is rational. We have $\sqrt{2} = \frac{q-p}{r-p}$. Since p, q, r are all rational (and $r - p \neq 0$), RHS is rational. This implies that LHS is rational, i.e. $\sqrt{2}$ is rational, which is a contradiction.

□

Example 1.9. Prove that every integer greater than 1 is divisible by a prime.

Proof. We proceed by a non-constructive proof.

If n is prime, then we are done as $n \mid n$.

If n is not prime, then n is composite. So n has a divisor d_1 such that $1 < d_1 < n$. If d_1 is prime then we are done as $d_1 \mid n$. If d_1 is not prime then d_1 is composite, has divisor d_2 such that $1 < d_2 < n$.

If d_2 is prime, then we are done as $d_2 \mid d_1$ and $d_1 \mid n$ imply $d_2 \mid n$. If d_2 is not prime then d_2 is composite, has divisor d_3 such that $1 < d_3 < d_2$.

Continuing in this manner after k times, we will get

$$1 < d_k < d_{k-1} < \cdots < d_2 < d_1 < n$$

where $d_i \mid n$ for all i .

Since there can only be a finite number of d_i 's between 1 and n , this process must stop after finite steps. On the other hand, the process will stop only if there is a d_i which is a prime. Hence we conclude that there must be a divisor d_i of n that is prime. □

Remark. This proof is also known as *proof by infinite descent*, a method which relies on the well-ordering principle on \mathbb{N} .

Example 1.10. Prove that the equation $x^2 + y^2 = 3z^2$ has no solutions (x, y, z) in integers where $z \neq 0$.

Proof. Suppose (x, y, z) is a solution. WLOG assume $z > 0$. By the least integer principle, we may also assume that our solution has z minimal. Taking remainders modulo 3, we see that

$$x^2 + y^2 \equiv 0 \pmod{3}$$

Since perfect squares can only be congruent to 0 or 1 modulo 3, we must have $x \equiv y \equiv 0 \pmod{3}$. Writing $x = 3a$ and $y = 3b$ for $a, b \in \mathbb{Z}$ gives

$$9a^2 + 9b^2 = 3z^2 \implies 3(a^2 + b^2) = z^2 \implies 3 \mid z^2 \implies 3 \mid z$$

Now let $z = 3c$ and cancel 3's to obtain

$$a^2 + b^2 = 3c^2.$$

We have therefore constructed another solution $(a, b, c) = (\frac{x}{3}, \frac{y}{3}, \frac{z}{3})$, but $0 < c < z$ contradicts the minimality of z . \square

Proof by Mathematical Induction

Induction is an extremely powerful method of proof used throughout mathematics. It deals with infinite families of statements which come in the form of lists. The idea behind induction is in showing how each statement follows from the previous one on the list—all that remains is to kick off this logical chain reaction from some starting point.

The *well-ordering principle* on \mathbb{N} states the following: every non-empty subset $S \subset \mathbb{N}$ has a smallest element; that is, there exists $m \in S$ such that $m \leq k$ for all $k \in S$.

The *principle of induction* states the following: Let $S \subset \mathbb{N}$. If (i) $1 \in S$, and (ii) $k \in S \implies k + 1 \in S$, then $S = \mathbb{N}$.

Lemma 1.11. *The well-ordering principle is equivalent to the principle of induction.*

Proof.

\implies Suppose otherwise, for a contradiction, that S exists with the given properties in the principle of induction, but $S \neq \mathbb{N}$.

Consider the set $\mathbb{N} \setminus S$. Then $\mathbb{N} \setminus S$ is not empty. By the well-ordering principle, $\mathbb{N} \setminus S$ has a least element p . Since $1 \in S$, $1 \notin \mathbb{N} \setminus S$ so $p \neq 1$, thus we must have $p > 1$.

Now consider $p - 1$. Since p is the least element of $\mathbb{N} \setminus S$, $p - 1 \notin \mathbb{N} \setminus S$ so $p - 1 \in S$. But by (ii) of the principle of induction, $p - 1 \in S$ implies $p \in S$, which contradicts the fact that $p \in \mathbb{N} \setminus S$.

\impliedby Suppose the principle of induction is true. Then this implies that Theorem 1.12 is true, which in turn implies that Theorem 1.16 is true. In order to prove the well-ordering of \mathbb{N} , we prove the following statement $P(n)$ by strong induction on n : If $S \subset \mathbb{N}$ and $n \in S$, then S has a least element.

The basis step is true, because if $1 \in S$ then 1 is the smallest element of S , since there are no smaller elements of \mathbb{N} .

Now suppose that $P(k)$ is true for $k = 1, \dots, n$. To show that $P(n + 1)$ is true, let $S \subset \mathbb{N}$ contain $n + 1$. If $n + 1$ is the smallest element of S , then we are done. Otherwise, S has a smaller element k , and $P(k)$ is true by the inductive hypothesis, so again S has a smallest element.

Hence by strong induction, $P(n)$ is true for all $n \in \mathbb{N}$. This implies the well-ordering of \mathbb{N} , because if S is a non-empty subset of \mathbb{N} , then pick $n \in S$. Since $n \in \mathbb{N}$, $P(n)$ is true, and therefore S has a smallest element. \square

Theorem 1.12 (Principle of mathematical induction). *Let $P(n)$ be a family of statements indexed by \mathbb{N} . Suppose that*

(i) $P(1)$ is true;

(ii) for all $k \in \mathbb{N}$, $P(k) \implies P(k + 1)$.

Then $P(n)$ is true for all $n \in \mathbb{N}$.

(i) is known as the *base case*; (ii) is known as the *inductive step*, where we assume $P(k)$ to be true—this is called the *inductive hypothesis*—and show that $P(k + 1)$ is true.

Proof. Apply the principle of induction to the set $S = \{n \in \mathbb{N} \mid P(n) \text{ is true}\}$. \square

We illustrate the application of this proving technique using a classic example.

Example 1.13. Prove that for any $n \in \mathbb{N}$,

$$\sum_{i=1}^n i = \frac{n(n+1)}{2}.$$

Proof. Induct on n . Let $P(n) : \sum_{i=1}^n i = \frac{n(n+1)}{2}$.

Clearly $P(1)$ holds. Now suppose $P(k)$ holds for some $k \in \mathbb{N}$, $k \geq 1$; that is,

$$\sum_{i=1}^k i = \frac{k(k+1)}{2}.$$

Adding $k+1$ to both sides,

$$\begin{aligned} \sum_{i=1}^{k+1} i &= \frac{k(k+1)}{2} + (k+1) \\ &= \frac{(k+1)(k+2)}{2} \\ &= \frac{(k+1)[(k+1)+1]}{2} \end{aligned}$$

thus $P(k+1)$ is true. Hence by induction, the result holds. \square

Example 1.14 (Bernoulli's inequality). Let $x \in \mathbb{R}$, $x > -1$. Then for all $n \in \mathbb{N}$,

$$(1+x)^n \geq 1+nx.$$

Proof. Induct on n . Fix $x > -1$. Let $P(n) : (1+x)^n \geq 1+nx$.

The base case $P(1)$ is clear. Suppose that $P(k)$ is true for some $k \in \mathbb{Z}^+$, $k \geq 1$. That is, $(1+x)^k \geq 1+kx$. Note that $1+x > 0$, and $kx^2 \geq 0$ (since $k > 0$ and $x^2 \geq 0$). Then

$$\begin{aligned} (1+x)^{k+1} &= (1+x)(1+x)^k \\ &\geq (1+x)(1+kx) && \text{[induction hypothesis]} \\ &= 1+(k+1)x+kx^2 \\ &\geq 1+(k+1)x && [\cdot \cdot kx^2 \geq 0] \end{aligned}$$

so $P(k+1)$ is true. Hence by induction, the result holds. \square

A corollary of induction is if the family of statements holds for $n \geq N$, rather than necessarily $n \geq 0$:

Corollary 1.15. Let $P(n)$ be a family of statements indexed by integers $n \geq N$ for some $N \in \mathbb{Z}$.

Suppose that

(i) $P(N)$ is true;

(ii) for all $k \geq N$, $P(k) \implies P(k+1)$.

Then $P(n)$ is true for all $n \geq N$.

Proof. Apply Theorem 1.12 to the statement $Q(n) = P(n + N)$ for $n \in \mathbb{N}$. \square

Another variant on induction is when the inductive step relies on some earlier case(s) but not necessarily the immediately previous case.

Theorem 1.16 (Strong induction). *Let $P(n)$ be a family of statements indexed by \mathbb{N} . Suppose that*

(i) $P(1)$ is true;

(ii) for all $k \in \mathbb{N}$, $P(1) \wedge \cdots \wedge P(k) \implies P(k + 1)$.

Then $P(n)$ is true for all $n \in \mathbb{N}$.

Proof. Let $Q(n)$ be the statement “ $P(k)$ holds for $k = 1, \dots, n$ ”. Then the conditions for the strong form are equivalent to (i) $Q(1)$ holds and (ii) for $n \in \mathbb{N}$, $Q(n) \implies Q(n + 1)$. By Theorem 1.12, $Q(n)$ holds for all $n \in \mathbb{N}$, and hence $P(n)$ holds for all n . \square

Example 1.17 (Fundamental theorem of arithmetic). Prove that every natural number greater than 1 may be expressed as a product of one or more prime numbers.

Proof. Let $P(n)$ be the statement that n may be expressed as a product of prime numbers.

Clearly $P(2)$ holds, since 2 is itself prime. Let $n \geq 2$ be a natural number and suppose that $P(k)$ holds for all $k < n$.

- If n is prime then it is trivially the product of the single prime number n .
- If n is not prime, then there must exist some $r, s > 1$ such that $n = rs$. By the inductive hypothesis, each of r and s can be written as a product of primes, and therefore $n = rs$ is also a product of primes.

In both cases, $P(n)$ holds. Hence by strong induction, $P(n)$ is true for all $n \in \mathbb{N}$. \square

The following is also another variant on induction.

Theorem 1.18 (Cauchy induction). *Let $P(n)$ be a family of statements indexed by $\mathbb{N}_{\geq 2}$. Suppose that*

(i) $P(2)$ is true;

(ii) for all $k \in \mathbb{N}$, $P(k) \implies P(2k)$ and $P(k) \implies P(k - 1)$.

Then $P(n)$ is true for all $n \in \mathbb{N}_{\geq 2}$.

Example 1.19 (AM–GM inequality). Given $n \in \mathbb{N}$, prove that for positive reals a_1, a_2, \dots, a_n ,

$$\frac{a_1 + a_2 + \cdots + a_n}{n} \geq \sqrt[n]{a_1 a_2 \cdots a_n}.$$

Proof. Let $P(n) : \frac{a_1 + a_2 + \cdots + a_n}{n} \geq \sqrt[n]{a_1 a_2 \cdots a_n}$.

Base case $P(2)$ is true because

$$\frac{a_1 + a_2}{2} \geq \sqrt{a_1 a_2} \iff (a_1 + a_2)^2 \geq 4a_1 a_2 \iff (a_1 - a_2)^2 \geq 0$$

Next we show that $P(n) \implies P(2n)$

$$\begin{aligned} \frac{a_1 + a_2 + \cdots + a_{2n}}{2n} &= \frac{\frac{a_1 + a_2 + \cdots + a_n}{n} + \frac{a_{n+1} + a_{n+2} + \cdots + a_{2n}}{n}}{2} \\ \frac{\frac{a_1 + a_2 + \cdots + a_n}{n} + \frac{a_{n+1} + a_{n+2} + \cdots + a_{2n}}{n}}{2} &\geq \frac{\sqrt[n]{a_1 a_2 \cdots a_n} + \sqrt[n]{a_{n+1} a_{n+2} \cdots a_{2n}}}{2} \\ \frac{\sqrt[n]{a_1 a_2 \cdots a_n} + \sqrt[n]{a_{n+1} a_{n+2} \cdots a_{2n}}}{2} &\geq \sqrt{\sqrt[n]{a_1 a_2 \cdots a_n} \sqrt[n]{a_{n+1} a_{n+2} \cdots a_{2n}}} \\ &= \sqrt{\sqrt[n]{a_1 a_2 \cdots a_n} \sqrt[n]{a_{n+1} a_{n+2} \cdots a_{2n}}} = \sqrt[2n]{a_1 a_2 \cdots a_{2n}} \end{aligned}$$

The first inequality follows from n -variable AM–GM, which is true by assumption, and the second inequality follows from 2-variable AM–GM, which is proven above.

Finally we show that $P(n) \implies P(n-1)$. By n -variable AM–GM, $\frac{a_1 + a_2 + \cdots + a_n}{n} \geq \sqrt[n]{a_1 a_2 \cdots a_n}$. Let $a_n = \frac{a_1 + a_2 + \cdots + a_{n-1}}{n-1}$. Then we have

$$\frac{a_1 + a_2 + \cdots + a_{n-1} + \frac{a_1 + a_2 + \cdots + a_{n-1}}{n-1}}{n} = \frac{a_1 + a_2 + \cdots + a_{n-1}}{n-1}$$

So,

$$\begin{aligned} \frac{a_1 + a_2 + \cdots + a_{n-1}}{n-1} &\geq \sqrt[n]{a_1 a_2 \cdots a_{n-1} \cdot \frac{a_1 + a_2 + \cdots + a_{n-1}}{n-1}} \\ \Rightarrow \left(\frac{a_1 + a_2 + \cdots + a_{n-1}}{n-1} \right)^n &\geq a_1 a_2 \cdots a_{n-1} \cdot \frac{a_1 + a_2 + \cdots + a_{n-1}}{n-1} \\ &\Rightarrow \left(\frac{a_1 + a_2 + \cdots + a_{n-1}}{n-1} \right)^{n-1} \geq a_1 a_2 \cdots a_{n-1} \\ &\Rightarrow \frac{a_1 + a_2 + \cdots + a_{n-1}}{n-1} \geq \sqrt[n-1]{a_1 a_2 \cdots a_{n-1}} \end{aligned}$$

By Cauchy induction, this proves the AM–GM inequality for n variables. □

Pigeonhole Principle

Theorem 1.20 (Pigeonhole principle). *If $kn + 1$ objects are distributed among n boxes, one of the boxes will contain at least $k + 1$ objects.*

Example 1.21 (IMO 1972). Prove that every set of 10 two-digit integer numbers has two disjoint subsets with the same sum of elements.

Proof. Let S be the set of 10 numbers. It has $2^{10} - 2 = 1022$ subsets that differ from both S and the empty set. They are the “pigeons”.

If $A \subset S$, the sum of elements of A cannot exceed $91 + 92 + \cdots + 99 = 855$. The numbers between 1 and 855, which are all possible sums, are the “holes”.

Because the number of “pigeons” exceeds the number of “holes”, there will be two “pigeons” in the same “hole”. Specifically, there will be two subsets with the same sum of elements. Deleting the common elements, we obtain two disjoint sets with the same sum of elements. \square

Example 1.22 (Putnam 2006). Prove that for every set $X = \{x_1, x_2, \dots, x_n\}$ of n real numbers, there exists a nonempty subset S of X and an integer m such that

$$\left| m + \sum_{x \in S} s \right| \leq \frac{1}{n+1}.$$

Proof. Recall that the fractional part of a real number x is $x - \lfloor x \rfloor$. Consider the fractional parts of the numbers $x_1, x_1 + x_2, \dots, x_1 + x_2 + \cdots + x_n$.

- If any of them is either in the interval $\left[0, \frac{1}{n+1}\right]$ or $\left[\frac{n}{n+1}, 1\right]$, then we are done.
- If not, consider these n numbers as the “pigeons” and the $n - 1$ intervals

$$\left[\frac{1}{n+1}, \frac{2}{n+1}\right], \left[\frac{2}{n+1}, \frac{3}{n+1}\right], \dots, \left[\frac{n-1}{n+1}, \frac{n}{n+1}\right]$$

as the “holes”. By the pigeonhole principle, two of these sums, say $x_1 + x_2 + \cdots + x_k$ and $x_1 + x_2 + \cdots + x_{k+m}$, belong to the same interval. But then their difference $x_{k+1} + \cdots + x_{k+m}$ lies within a distance of $\frac{1}{n+1}$ of an integer, and we are done. \square

Exercises

Exercise 1.1. Negate the statement

for all real numbers x , if $x > 2$, then $x^2 > 4$

Solution. In logical notation, this statement is $(\forall x \in \mathbb{R})[x > 2 \implies x^2 > 4]$.

$$\begin{aligned} & \neg\{(\forall x \in \mathbb{R})[x > 2 \implies x^2 > 4]\} \\ & \equiv (\exists x \in \mathbb{R})\neg[x > 2 \implies x^2 > 4] \\ & \equiv (\exists x \in \mathbb{R})\neg[(x > 2) \vee (x^2 > 4)] \\ & \equiv (\exists x \in \mathbb{R})[(x > 2) \wedge (x^2 \leq 4)] \end{aligned}$$

□

Exercise 1.2. Negate surjectivity.

Solution. If $f : X \rightarrow Y$ is not surjective, then it means that there exists $y \in Y$ not in the image of X , i.e. for all x in X we have $f(x) \neq y$.

$$\begin{aligned} & \neg\forall y \in Y, \exists x \in X, f(x) = y \\ & \equiv \exists y \in Y, \neg(\exists x \in X, f(x) = y) \\ & \equiv \exists y \in Y, \forall x \in X, \neg(f(x) = y) \\ & \equiv \exists y \in Y, \forall x \in X, f(x) \neq y \end{aligned}$$

□

Exercise 1.3. Use the Unique Factorisation Theorem to prove that, if a positive integer n is not a perfect square, then \sqrt{n} is irrational.

[The Unique Factorisation Theorem states that every integer $n > 1$ has a unique standard factored form, i.e. there is exactly one way to express $n = p_1^{k_1} p_2^{k_2} \cdots p_t^{k_t}$ where $p_1 < p_2 < \cdots < p_t$ are distinct primes and k_1, k_2, \dots, k_t are some positive integers.]

Solution. Prove by contradiction. Suppose n is not a perfect square and \sqrt{n} is rational. Then $\sqrt{n} = \frac{a}{b}$ for some $a, b \in \mathbb{Z}$. Squaring both sides and clearing denominator gives

$$nb^2 = a^2. \quad (*)$$

Consider the standard factored forms of n , a and b :

$$\begin{aligned} n &= p_1^{k_1} p_2^{k_2} \cdots p_t^{k_t} \\ a &= q_1^{e_1} q_2^{e_2} \cdots q_u^{e_u} \implies a^2 = q_1^{2e_1} q_2^{2e_2} \cdots q_u^{2e_u} \\ b &= r_1^{f_1} r_2^{f_2} \cdots r_v^{f_v} \implies b^2 = r_1^{2f_1} r_2^{2f_2} \cdots r_v^{2f_v} \end{aligned}$$

i.e. the powers of primes in the standard factored form of a^2 and b^2 are all even integers.

This means the powers k_i of primes p_i in the standard factored form of n are also even by Unique Factorisation Theorem. Note that all p_i appear in the standard factored form of a^2 with even power $2c_i$, because of (*). By UFT, p_i must also appear in the standard factored form of nb^2 with the same even power $2c_i$.

If $p_i \nmid b$, then $k_i = 2c_i$ which is even. If $p_i \mid b$, then p_i will appear in b^2 with even power $2d_i$. So $k_i + 2d_i = 2c_i$, and hence $k_i = 2(c_i - d_i)$, which is again even.

$$\text{Hence } n = p_1^{k_1} p_2^{k_2} \cdots p_t^{k_t} = \left(p_1^{\frac{k_1}{2}} p_2^{\frac{k_2}{2}} \cdots p_t^{\frac{k_t}{2}} \right)^2.$$

Since $\frac{k_i}{2}$ are all integers, $p_1^{\frac{k_1}{2}} p_2^{\frac{k_2}{2}} \cdots p_t^{\frac{k_t}{2}}$ is an integer and n is a perfect square. This contradicts the given hypothesis that n is not a perfect square. \square

Exercise 1.4. Prove that for every pair of irrational numbers p and q such that $p < q$, there is an irrational x such that $p < x < q$.

Solution. Consider the average of p and q , i.e., $\frac{p+q}{2}$. Evidently $p < \frac{p+q}{2} < q$.

Since it may not always be the case that $\frac{p+q}{2}$ is irrational (so we cannot immediately take $x = \frac{p+q}{2}$), we need to consider two cases:

$\frac{p+q}{2}$ **is irrational** Take $x = \frac{p+q}{2}$ and we are done.

$\frac{p+q}{2}$ **is rational** Let $r = \frac{p+q}{2}$, and take the average of p and r , i.e., $\frac{p+r}{2}$. Evidently $p < \frac{p+r}{2} < r < q$.

Since p is irrational and r is rational, $\frac{p+r}{2}$ is irrational. In this case, take $x = \frac{3p+q}{4}$.

\square

Exercise 1.5. Given n real numbers a_1, a_2, \dots, a_n . Show that there exists an a_i ($1 \leq i \leq n$) such that a_i is greater than or equal to the mean of the n numbers.

Solution. Prove by contradiction.

Let \bar{a} denote the mean value of the n given numbers. Suppose $a_i < \bar{a}$ for all a_i . Then

$$\bar{a} = \frac{a_1 + a_2 + \cdots + a_n}{n} < \frac{\bar{a} + \bar{a} + \cdots + \bar{a}}{n} = \frac{n\bar{a}}{n} = \bar{a}.$$

We derive $\bar{a} < \bar{a}$, which is a contradiction.

Hence there must be some a_i such that $a_i \geq \bar{a}$. \square

Exercise 1.6. Prove that the following statement is false: there is an irrational number a such that for all irrational number b , ab is rational.

Idea. Prove the negation of the statement: for every irrational number a , there is an irrational number b such that ab is irrational. We shall adopt a constructive proof (note that we can consider multiple cases and construct more than one b).

Solution. Given an irrational number a , let us consider $\frac{\sqrt{2}}{a}$. We consider cases:

- If $\frac{\sqrt{2}}{a}$ is irrational, take $b = \frac{\sqrt{2}}{a}$. Then $ab = \sqrt{2}$ which is irrational.

- If $\frac{\sqrt{2}}{a}$ is rational, its reciprocal $\frac{a}{\sqrt{2}}$ is rational. Since $\sqrt{6}$ is irrational, the product $\left(\frac{a}{\sqrt{2}}\right)\sqrt{6} = a\sqrt{3}$ is irrational. Take $b = \sqrt{3}$, which is irrational. Then $ab = a\sqrt{3}$ is irrational.

□

Exercise 1.7. Prove that there are infinitely many prime numbers that are congruent to 3 modulo 4.

Idea. It is not really possible to come up with a direct proof, so we prove by contradiction.

Solution. Suppose, for a contradiction, that there are only finitely many primes that are congruent to 3 modulo 4. Let p_1, p_2, \dots, p_m be the list of all the primes that are congruent to 3 modulo 4.

Let $M = (p_1 p_2 \cdots p_m)^2 + 2$.

We have the following observation:

- (i) $M \equiv 3 \pmod{4}$.
- (ii) Every p_i divides $M - 2$.
- (iii) None of the p_i divides M . [Otherwise, together with (ii), this will imply p_i divides 2, which is impossible.]
- (iv) M is not a prime number. [Otherwise, by (i), M is a prime number congruent to 3 modulo 4. But $M \neq p_i$ for all $1 \leq i \leq m$. This contradicts the assumption that p_1, p_2, \dots, p_m are all the prime numbers congruent to 3 modulo 4.]

From the above discussion, we know that M is a composite number by (iv). So it has a prime factorization $M = q_1 q_2 \cdots q_k$.

Since M is odd, all these prime factors q_j must be odd, and hence q_j must be congruent to either 1 or 3 modulo 4.

By (iii), q_j cannot be any of the p_i . So all q_j must be congruent to 1 modulo 4. Then M , which is the product of q_j , must also be congruent to 1 modulo 4.

This contradicts (i) that M is congruent to 3 modulo 4.

Hence we conclude that there must be infinitely many primes that are congruent to 3 modulo 4. □

Exercise 1.8. Prove that, for any positive integer n , there exists a perfect square m^2 such that $n \leq m^2 \leq 2n$.

Idea. A direct proof by construction is not quite possible, so we prove by contradiction.

Solution. Suppose, for a contradiction, that $n > m^2$ and $(m + 1)^2 > 2n$ for some positive integer n , so that there is no square between n and $2n$. Then

$$(m + 1)^2 > 2n > 2m^2.$$

Since we are dealing with integers and the inequalities are strict, we get

$$(m + 1)^2 \geq 2m^2 + 2$$

which simplifies to

$$0 \geq m^2 - 2m + 1 = (m - 1)^2$$

The only value for which this is possible is $m = 1$, but you can eliminate that easily enough. □

Exercise 1.9. Prove that for every positive integer $n \geq 4$,

$$n! > 2^n.$$

Solution. Induct on n . Let $P(n) : n! > 2^n$.

The base case $P(4)$ is clear. Now suppose $P(k)$ is true for some $k \in \mathbb{N}_{\geq 4}$, i.e., $k! > 2^k$. Then

$$(k+1)! = k!(k+1) > 2^k(k+1) > 2^k \cdot 2 = 2^{k+1},$$

so $P(k+1)$ is true. □

Exercise 1.10. Prove by mathematical induction, for $n \geq 2$,

$$\sqrt[n]{n} < 2 - \frac{1}{n}.$$

Solution. Induct on n . Let $P(n) : \sqrt[n]{n} < 2 - \frac{1}{n}$, for $n \geq 2$.

The base case $P(2)$ is clear. Now assume $P(k)$ is true for $k \geq 2$, $k \in \mathbb{N}$, i.e., $\sqrt[k]{k} < 2 - \frac{1}{k}$, or

$$k < \left(2 - \frac{1}{k}\right)^k.$$

We want to prove that $P(k+1)$ is true; that is,

$$k+1 < \left(2 - \frac{1}{k+1}\right)^{k+1}$$

Since $k > 2$, we have

$$\begin{aligned} \left(2 - \frac{1}{k+1}\right)^{k+1} &> \left(2 - \frac{1}{k}\right)^{k+1} && [\because k > 2] \\ &= \left(2 - \frac{1}{k}\right)^k \left(2 - \frac{1}{k}\right) \\ &> k \left(2 - \frac{1}{k}\right) && [\text{by inductive hypothesis}] \\ &= 2k - 1 > k - 1 \end{aligned}$$

so $P(k+1)$ is true. □

Exercise 1.11. Prove that, for all integers $n \geq 3$,

$$\left(1 + \frac{1}{n}\right)^n < n.$$

Solution. For the base case $P(3)$, $\left(1 + \frac{1}{3}\right)^3 = \frac{64}{27} = 2\frac{10}{27} < 3$. Hence $P(3)$ is true.

Assume that $P(k)$ is true for some $k \in \mathbb{N}_{\geq 3}$; that is,

$$\left(1 + \frac{1}{k}\right)^k < k.$$

Multiplying both sides by $\left(1 + \frac{1}{k}\right)$ (to get a $k + 1$ in the power),

$$\left(1 + \frac{1}{k}\right)^k \left(1 + \frac{1}{k}\right) = \left(1 + \frac{1}{k}\right)^{k+1} < k \left(1 + \frac{1}{k}\right) = k + 1$$

Since $k < k + 1 \iff \frac{1}{k} > \frac{1}{k + 1}$,

$$\left(1 + \frac{1}{k}\right)^{k+1} > \left(1 + \frac{1}{k+1}\right)^{k+1}$$

The rest of the proof follows easily. \square

A sequence of integers F_i , where integer $1 \leq i \leq n$, is called the *Fibonacci sequence* if and only if it is defined recursively by $F_1 = 1$, $F_2 = 1$, $F_n = F_{n-1} + F_{n-2}$ for $n > 2$.

Exercise 1.12. Let (a_n) be a sequence of integers defined recursively by the initial conditions $a_1 = 1$, $a_2 = 1$, $a_3 = 3$ and the recurrence relation $a_n = a_{n-1} + a_{n-2} + a_{n-3}$ for $n > 3$.

For all $n \in \mathbb{N}$, prove that

$$a_n \leq 2^{n-1}.$$

Idea. Given the recurrence relation, we may need to use *strong induction*: use $P(k)$, $P(k + 1)$, $P(k + 2)$ to prove $P(k + 3)$, for all $k \in \mathbb{N}$.

Solution. Let $P(n) : a_n \leq 2^{n-1}$.

The base cases $P(1)$, $P(2)$, $P(3)$ are clear. Now assume $P(k)$, $P(k + 1)$, $P(k + 2)$ are true, for some $k \in \mathbb{N}$. We will show that $P(k + 3)$ is true.

By the inductive hypothesis, for $k \in \mathbb{N}$ we have

$$a_k \leq 2^k, \quad a_{k+1} \leq 2^{k+1}, \quad a_{k+2} \leq 2^{k+2}.$$

Then

$$\begin{aligned} a_{k+3} &= a_k + a_{k+1} + a_{k+2} && \text{[start from recurrence relation]} \\ &\leq 2^k + 2^{k+1} + 2^{k+2} && \text{[use inductive hypothesis]} \\ &= 2^k(1 + 2 + 2^2) \\ &< 2^k(2^3) && \text{[approximation, since } 1 + 2 + 2^2 < 2^3\text{]} \\ &= 2^{k+3} \end{aligned}$$

which is precisely $P(k + 3) : a_{k+3} \leq 2^{k+3}$. \square

Exercise 1.13. For $m, n \in \mathbb{N}$, prove that

$$F_{n+m+1} = F_n F_m + F_{n+1} F_{m+1}.$$

Solution. Induct on n . Let $P(n) : F_{n+m+1} = F_n F_m + F_{n+1} F_{m+1}$ for all $m \in \mathbb{N}$ in the cases $k = n$ and $k = n + 1$.

To show that $P(0)$ is true, note that

$$F_{m+1} = F_0F_m + F_1F_{m+1}$$

and

$$F_{m+2} = F_1F_m + F_2F_{m+1}$$

for all m , as $F_0 = 0$ and $F_1 = F_2 = 1$.

Now assume $P(n)$ is true; that is, for all $m \in \mathbb{N}$,

$$\begin{aligned} F_{n+m+1} &= F_nF_m + F_{n+1}F_{m+1}, \\ F_{n+m+2} &= F_{n+1}F_m + F_{n+2}F_{m+1}. \end{aligned}$$

Then

$$\begin{aligned} F_{n+m+3} &= F_{n+m+2} + F_{n+m+1} \\ &= F_nF_m + F_{n+1}F_{m+1} + F_{n+1}F_m + F_{n+2}F_{m+1} \\ &= (F_n + F_{n+1})F_m + (F_{n+1} + F_{n+2})F_{m+1} \\ &= F_{n+2}F_m + F_{n+3}F_{m+1} \end{aligned}$$

thus $P(n+1)$ is true, for all $m \in \mathbb{N}$. □

2 Set Theory

Summary

- Basic definitions relating to sets (excluding detailed axiomatic discussions).
- Relations and related concepts including binary relation, partial order, total order, well order, equivalence relations, equivalence relations, equivalence class, quotient set, partition.
- Functions, injectivity, surjectivity, bijectivity, composition, invertibility.

§2.1 Basics of Naive Set Theory

Definitions and Notations

A **set** S can be loosely defined as a collection of objects¹. For a set S , we write $x \in S$ to mean that x is an **element** of S , and $x \notin S$ if otherwise.

To describe a set, one can list its elements explicitly. A set can also be defined in terms of some property $P(x)$ that the elements $x \in S$ satisfy, denoted by the following set builder notation:

$$\{x \in S \mid P(x)\}$$

The following sets of numbers are frequently encountered.

- $\mathbb{N} = \{1, 2, 3, \dots\}$ denotes the natural numbers (non-negative integers).
- $\mathbb{Z} = \{\dots, -2, -1, 0, 1, 2, \dots\}$ denotes the integers.
- $\mathbb{Q} = \left\{\frac{p}{q} \mid p, q \in \mathbb{Z}, q \neq 0\right\}$ denotes the rational numbers.
- \mathbb{R} denotes the real numbers (the construction of which, using Dedekind cuts, will be discussed in Chapter 10).
- $\mathbb{C} = \{x + yi \mid x, y \in \mathbb{R}\}$ denotes the complex numbers.

We have that $\mathbb{N} \subset \mathbb{Z} \subset \mathbb{Q} \subset \mathbb{R} \subset \mathbb{C}$.

The **empty set** is the set with no elements, denoted by \emptyset .

¹*Russell's paradox*, after the mathematician and philosopher Bertrand Russell (1872–1970), provides a warning as to the looseness of our definition of a set. Suppose H is the collection of sets that are not elements of themselves; that is,

$$H = \{S \mid S \notin S\}.$$

The problem arises when we ask the question of whether or not H is itself in H ? On one hand, if $H \notin H$ then H meets the precise criterion for being in H and so $H \in H$, a contradiction. On the other hand, if $H \in H$ then by the property required for this to be the case, $H \notin H$, another contradiction. Thus we have a paradox: H is neither in H , nor not in H .

The modern resolution of Russell's paradox is that we have taken too naive an understanding of "collection", and that Russell's "set" H is in fact not a set. It does not fit within axiomatic set theory (which relies on the so-called ZF axioms), and so the question of whether or not H is in H simply doesn't make sense.

A is a **subset** of B if every element of A is in B , denoted by $A \subset B$:

$$A \subset B \iff (\forall x)(x \in A \implies x \in B)$$

We denote $A \subsetneq B$ to explicitly mean that $A \subset B$ and $A \neq B$; we call A a *proper subset* of B .

Lemma 2.1 (\subset is transitive). *If $A \subset B$ and $B \subset C$, then $A \subset C$.*

Proof. Let $x \in A$. Since $A \subset B$ and $x \in A$, $x \in B$. Since $B \subset C$ and $x \in B$, $x \in C$. Hence $A \subset C$. \square

A and B are **equal** if and only if they contain the same elements, denoted by $A = B$.

Lemma 2.2 (Double inclusion). *Let $A \subset S$ and $B \subset S$. Then*

$$A = B \iff (A \subset B) \wedge (B \subset A)$$

Proof. We have

$$\begin{aligned} A = B &\iff (\forall x)[x \in A \iff x \in B] \\ &\iff (\forall x)[(x \in A \implies x \in B) \wedge (x \in B \implies x \in A)] \\ &\iff \{(\forall x)[x \in A \implies x \in B]\} \wedge \{(\forall x)[x \in B \implies x \in A]\} \\ &\iff (A \subset B) \wedge (B \subset A) \end{aligned}$$

\square

Remark. Double inclusion is a useful tool to prove that two sets are equal.

Some frequently occurring subsets of \mathbb{R} are known as **intervals**, which can be visualised as sections of the real line. We define *bounded intervals*

$$\begin{aligned} (a, b) &= \{x \in \mathbb{R} \mid a < x < b\}, \\ [a, b] &= \{x \in \mathbb{R} \mid a \leq x \leq b\}, \\ [a, b) &= \{x \in \mathbb{R} \mid a \leq x < b\}, \\ (a, b] &= \{x \in \mathbb{R} \mid a < x \leq b\}, \end{aligned}$$

and *unbounded intervals*

$$\begin{aligned} (a, \infty) &= \{x \in \mathbb{R} \mid a < x\}, \\ [a, \infty) &= \{x \in \mathbb{R} \mid a \leq x\}, \\ (-\infty, a) &= \{x \in \mathbb{R} \mid x < a\}, \\ (-\infty, a] &= \{x \in \mathbb{R} \mid x \leq a\}. \end{aligned}$$

An interval of the first type (a, b) is called an *open interval*; an interval of the second type $[a, b]$ is called a *closed interval*. Note that if $a = b$, then $[a, b] = \{a\}$, while $(a, b) = [a, b) = (a, b] = \emptyset$.

The **power set** $\mathcal{P}(A)$ of A is the set of all subsets of A (including the set itself and the empty set):

$$\mathcal{P}(A) = \{S \mid S \subset A\}.$$

An **ordered pair** is denoted by (a, b) , where the order of the elements matters. Two pairs (a_1, b_1) and (a_2, b_2) are equal if and only if $a_1 = a_2$ and $b_1 = b_2$. Similarly, we have ordered triples (a, b, c) , quadruples (a, b, c, d) and so on. If there are n elements it is called an n -tuple.

The **Cartesian product** of sets A and B is the set of all ordered pairs with the first element of the pair coming from A and the second from B :

$$A \times B := \{(a, b) \mid a \in A, b \in B\}.$$

More generally, we define $A_1 \times A_2 \times \cdots \times A_n$ to be the set of all ordered n -tuples (a_1, a_2, \dots, a_n) , where $a_i \in A_i$ for $1 \leq i \leq n$. If all the A_i are the same, we write the product as A^n .

Example 2.3. \mathbb{R}^2 is the Euclidean plane, \mathbb{R}^3 is the Euclidean space, and \mathbb{R}^n is the n -dimensional Euclidean space.

$$\begin{aligned}\mathbb{R} \times \mathbb{R} &= \mathbb{R}^2 = \{(x, y) \mid x, y \in \mathbb{R}\} \\ \mathbb{R} \times \mathbb{R} \times \mathbb{R} &= \mathbb{R}^3 = \{(x, y, z) \mid x, y, z \in \mathbb{R}\} \\ \mathbb{R}^n &= \{(x_1, x_2, \dots, x_n) \mid x_1, x_2, \dots, x_n \in \mathbb{R}\}\end{aligned}$$

Lemma 2.4. Let A, B, C, D be sets.

- (i) $A \times \emptyset = \emptyset \times A = \emptyset$.
- (ii) $A \times (B \cup C) = (A \times B) \cup (A \times C)$.
- (iii) $A \times (B \cap C) = (A \times B) \cap (A \times C)$.
- (iv) $(A \cap B) \times (C \cap D) = (A \times C) \cap (B \times D)$.
- (v) $(A \cup B) \times (C \cup D) \subset (A \times C) \cup (B \times D)$.

Proof.

- (i) Evidently $\emptyset \subset A \times \emptyset$, which is vacuously true.

To show the opposite containment,

- (ii)
- (iii)
- (iv)
- (v)

□

Algebra of Sets

We now discuss the algebra of sets. Given $A \subset S$ and $B \subset S$,

- (i) The **union** $A \cup B$ is the set consisting of elements that are in A or B (or both):

$$A \cup B = \{x \in S \mid x \in A \vee x \in B\}$$

(ii) The **intersection** $A \cap B$ is the set consisting of elements that are in both A and B :

$$A \cap B = \{x \in S \mid x \in A \wedge x \in B\}$$

A and B are **disjoint** if both sets have no element in common: $A \cap B = \emptyset$.

More generally, we can take unions and intersections of arbitrary numbers of sets (could be finitely or infinitely many). Given a family of sets $\{A_i \mid i \in I\}$ where I is an *indexing set*, we write

$$\bigcup_{i \in I} A_i = \{x \mid \exists i \in I, x \in A_i\},$$

and

$$\bigcap_{i \in I} A_i = \{x \mid \forall i \in I, x \in A_i\}.$$

(iii) The **complement** of A , denoted by A^c , is the set containing elements that are not in A :

$$A^c = \{x \in S \mid x \notin A\}$$

(iv) The **set difference**, or complement of B in A , denoted by $A \setminus B$, is the subset consisting of those elements that are in A and not in B :

$$A \setminus B = \{x \in A \mid x \notin B\}$$

Note that $A \setminus B = A \cap B^c$.

Lemma 2.5 (Distributive laws). *Let $A, B, C \subset S$. Then*

$$(i) \quad A \cup (B \cap C) = (A \cup B) \cap (A \cup C);$$

$$(ii) \quad A \cap (B \cup C) = (A \cap B) \cup (A \cap C).$$

Proof.

(i) Suppose $x \in A \cup (B \cap C)$. Then

$$\begin{aligned} x \in A \cup (B \cap C) &\iff x \in A \quad \vee \quad x \in B \cap C \\ &\iff x \in A \quad \vee \quad (x \in B) \wedge (x \in C) \\ &\iff (x \in A) \vee (x \in B) \quad \wedge \quad (x \in A) \vee (x \in C) \\ &\iff x \in A \cup B \quad \wedge \quad x \in A \cup C \\ &\iff x \in (A \cup B) \cap (A \cup C). \end{aligned}$$

Thus $A \cup (B \cap C) \subset (A \cup B) \cap (A \cup C)$.

Conversely suppose that $x \in (A \cup B) \cap (A \cup C)$. Then go in the reverse direction of the above steps to show that $(A \cup B) \cap (A \cup C) \subset A \cup (B \cap C)$.

By double inclusion, $(A \cup B) \cap (A \cup C) = A \cup (B \cap C)$.

(ii) Similar.

□

Lemma 2.6 (de Morgan's laws). *Let $A, B \subset S$. Then*

$$(i) \quad (A \cup B)^c = A^c \cap B^c;$$

$$(ii) \quad (A \cap B)^c = A^c \cup B^c.$$

Proof.

(i)

$$\begin{aligned} x \in (A \cup B)^c &\iff x \notin A \cup B \\ &\iff x \notin A \quad \wedge \quad x \notin B \\ &\iff x \in A^c \quad \wedge \quad x \in B^c \\ &\iff x \in A^c \cap B^c \end{aligned}$$

(ii) Similar.

□

De Morgan's laws extend naturally to any number of sets. Suppose $\{A_i \mid i \in I\}$ is a family of subsets of S , then

$$\begin{aligned} \left(\bigcap_{i \in I} A_i \right)^c &= \bigcup_{i \in I} A_i^c, \\ \left(\bigcup_{i \in I} A_i \right)^c &= \bigcap_{i \in I} A_i^c. \end{aligned}$$

Lemma 2.7. *The following hold:*

$$(i) \quad \left(\bigcup_{i \in I} A_i \right) \cup B = \bigcup_{i \in I} (A_i \cup B)$$

$$(ii) \quad \left(\bigcap_{i \in I} A_i \right) \cup B = \bigcap_{i \in I} (A_i \cup B)$$

$$(iii) \quad \left(\bigcup_{i \in I} A_i \right) \cup \left(\bigcup_{j \in J} B_j \right) = \bigcup_{(i,j) \in I \times J} (A_i \cup B_j)$$

$$(iv) \quad \left(\bigcap_{i \in I} A_i \right) \cup \left(\bigcap_{j \in J} B_j \right) = \bigcap_{(i,j) \in I \times J} (A_i \cup B_j)$$

§2.2 Relations

Definition and Examples

Definition 2.8 (Relation). R is a **relation** between A and B if $R \subset A \times B$; $a \in A$ and $b \in B$ are said to be *related* if $(a, b) \in R$, denoted aRb .

Remark. A relation is a set of ordered pairs.

Visually speaking, a relation is uniquely determined by a simple bipartite graph over A and B . On the bipartite graph, this is usually represented by an edge between a and b .

Example 2.9. In many cases we do not actually use R to write the relation because there is some other conventional notation:

- The “less than or equal to” relation \leq on the set of real numbers is

$$\{(x, y) \in \mathbb{R}^2 \mid x \leq y\} \subset \mathbb{R}^2;$$

we write $x \leq y$ if (x, y) is in this set.

- The “divides” relation $|$ on \mathbb{N} is

$$\{(m, n) \in \mathbb{N}^2 \mid m \text{ divides } n\} \subset \mathbb{N}^2;$$

we write $m \mid n$ if (m, n) is in this set.

- For a set S , the “subset” relation \subset on $\mathcal{P}(S)$ is

$$\{(A, B) \in \mathcal{P}(S)^2 \mid A \subset B\} \subset \mathcal{P}(S)^2;$$

we write $A \subset B$ if (A, B) is in this set.

If $A \times B$ is the smallest Cartesian product of which R is a subset, we call A and B the *domain* and *range* of R respectively, denoted by $\text{dom } R$ and $\text{ran } R$ respectively.

Example 2.10. Given $R = \{(1, a), (1, b), (2, b), (3, b)\}$, then $\text{dom } R = \{1, 2, 3\}$ and $\text{ran } R = \{a, b\}$.

Definition 2.11 (Binary relation). A **binary relation** in A is a relation between A and itself; that is, $R \subset A \times A$.

Properties of Relations

Let A be a set, R a relation on A , $x, y, z \in A$. We say that

- (i) R is **reflexive** if xRx for all $x \in A$;
- (ii) R is **symmetric** if $xRy \implies yRx$;

(iii) R is **anti-symmetric** if xRy and $yRx \implies x = y$;

(iv) R is **transitive** if xRy and $yRz \implies xRz$.

Example 2.12 (Less than or equal to). The relation \leq on R is reflexive, anti-symmetric, and transitive, but not symmetric.

Definition 2.13. A **partial order** on a non-empty set A is a relation on A satisfying reflexivity, anti-symmetry and transitivity.

A **total order** on A is a partial order on A such that if for every $x, y \in A$, either xRy or yRx .

A **well order** on A is a total order on A such that every non-empty subset of A has a minimal element; that is, for each non-empty $B \subset A$ there exists $s \in B$ such that $s \leq b$ for all $b \in B$.

Example 2.14.

- Less than: the relation $<$ on R is not reflexive, symmetric, or anti-symmetric, but it is transitive.
- Not equal to: the relation \neq on R is not reflexive, anti-symmetric or transitive, but it is symmetric.

Equivalence Relations

One important type of relation is an equivalence relation. An equivalence relation is a way of saying two objects are, in some particular sense, “the same”.

Definition 2.15 (Equivalence relation). A relation \sim on a set A is an **equivalence relation** if it is reflexive, symmetric and transitive.

Notation. We denote $a \sim b$ for $(a, b) \in R$.

An equivalence relation provides a way of grouping together elements that can be viewed as being the same:

Definition 2.16 (Equivalence class). Given an equivalence relation \sim on a set A , and given $x \in A$, the **equivalence class** of x is

$$[x] := \{y \in A \mid y \sim x\}.$$

Grouping the elements of a set into equivalence classes provides a partition of the set, which we define as follows:

Definition 2.17 (Partition). A **partition** of a set A is a collection of subsets $\{A_i \subset A \mid i \in I\}$, where I is an indexing set, with the property that

- $A_i \neq \emptyset$ for all $i \in I$ (all the subsets are non-empty)
- $\bigcup_{i \in I} A_i = A$ (every member of A lies in one of the subsets)
- $A_i \cap A_j = \emptyset$ for every $i \neq j$ (the subsets are disjoint)

The subsets are called the *parts* of the partition.

Proposition 2.18. *Let \sim be an equivalence relation on a non-empty set X . Then the equivalence classes under \sim are a partition of X .*

To prove this, we need to show that

- (i) every equivalence class is non-empty;
- (ii) every element of X is an element of an equivalence class;
- (iii) every element of X lies in exactly one equivalence class.

Proof.

- (i) An equivalence class $[x]$ contains x as $x \sim x$, by reflexivity of the relation. Thus $[x] \neq \emptyset$.
- (ii) From (i), note that every $x \in X$ is in the equivalence class $[x]$, so every element of X is an element of at least one equivalence class.
- (iii) Suppose otherwise, for a contradiction, that some element of X lies in more than one equivalence class. Let $x \in X$ such that $x \in [y]$ and $x \in [z]$; we want to show that $[y] = [z]$ (using double inclusion).
Let $a \in [y]$, so $a \sim y$. Also $x \in [y]$ so $x \sim y$. By symmetry, $y \sim x$. By transitivity, $a \sim x$. Now $x \in [z]$ so $x \sim z$ and similarly $a \sim z$ thus $a \in [z]$. Hence $[y] \subset [z]$.
By the same argument, $[z] \subset [y]$. Hence $[y] = [z]$.

□

Definition 2.19 (Quotient set). The *quotient set* is the set of all equivalence classes, denoted by A/\sim .

Example 2.20 (Modular arithmetic). Fix a positive integer n . Define a relation on \mathbb{Z} by

$$a \sim b \iff n \mid (b - a).$$

Lemma. \sim is a equivalence relation.

Proof. Let $a, b \in \mathbb{Z}$.

- (i) $a \sim a$ so \sim is reflexive.
- (ii) $a \sim b \implies b \sim a$ for any integers a and b , so \sim is symmetric.
- (iii) If $a \sim b$ and $b \sim c$, then $n \mid (a - b)$ and $n \mid (b - c)$, so $n \mid (a - b) + (b - c) = (a - c)$, so $a \sim c$ and \sim is transitive.

□

Notation. We write $a \equiv b \pmod{n}$ if $a \sim b$.

For any $k \in \mathbb{Z}$ we denote the equivalence class of a by $[a]$, called the *congruence class* (or *residue class*) of a mod n , which consists of the integers which differ from a by an integral multiple of n ; that is,

$$[a] = \{a + kn \mid k \in \mathbb{Z}\}.$$

There are precisely n distinct congruence classes mod n , namely

$$[0], [1], \dots, [n-1],$$

determined by the possible remainders after division by n ; and these residue classes partition the integers \mathbb{Z} . The set of equivalence classes under this equivalence relation is denoted by $\mathbb{Z}/n\mathbb{Z}$, and called the *integers modulo n* .

Define addition and multiplication on $\mathbb{Z}/n\mathbb{Z}$ as follows: for $[a], [b] \in \mathbb{Z}/n\mathbb{Z}$,

$$\begin{aligned} [a] + [b] &= [a + b] \\ [a][b] &= [ab]. \end{aligned}$$

This means that to compute the sum / product of two elements $[a], [b] \in \mathbb{Z}/n\mathbb{Z}$, take any *representative* $a \in [a]$, $b \in [b]$, and add / multiply integers a and b as usual in \mathbb{Z} , then take the congruence class containing the result.

Lemma. *Addition and multiplication on $\mathbb{Z}/n\mathbb{Z}$ are well-defined; that is, they do not depend on the choices of representatives for the classes involved. More precisely, if $a_1, a_2 \in \mathbb{Z}$ and $b_1, b_2 \in \mathbb{Z}$ with $\overline{a_1} = \overline{b_1}$ and $\overline{a_2} = \overline{b_2}$, then $\overline{a_1 + a_2} = \overline{b_1 + b_2}$ and $\overline{a_1 a_2} = \overline{b_1 b_2}$, i.e., If*

$$a_1 \equiv b_1 \pmod{n}, \quad a_2 \equiv b_2 \pmod{n}$$

then

$$a_1 + a_2 \equiv b_1 + b_2 \pmod{n}, \quad a_1 a_2 \equiv b_1 b_2 \pmod{n}.$$

Proof. Suppose $a_1 \equiv b_1 \pmod{n}$, i.e., $n \mid (a_1 - b_1)$. Then $a_1 = b_1 + sn$ for some integer s . Similarly, $a_2 \equiv b_2 \pmod{n}$ means $a_2 = b_2 + tn$ for some integer t .

Then $a_1 + a_2 = (b_1 + b_2) + (s + t)n$ so that $a_1 + a_2 \equiv b_1 + b_2 \pmod{n}$, which shows that the sum of the residue classes is independent of the representatives chosen.

Similarly, $a_1 a_2 = (b_1 + sn)(b_2 + tn) = b_1 b_2 + (b_1 t + b_2 s + stn)n$ shows that $a_1 a_2 \equiv b_1 b_2 \pmod{n}$ and so the product of the residue classes is also independent of the representatives chosen. \square

An important subset of $\mathbb{Z}/n\mathbb{Z}$ consists of the collection of congruence classes which have a multiplicative inverse in $\mathbb{Z}/n\mathbb{Z}$:

$$(\mathbb{Z}/n\mathbb{Z})^\times := \{[a] \in \mathbb{Z}/n\mathbb{Z} \mid \exists [c] \in \mathbb{Z}/n\mathbb{Z}, [a][c] = [1]\}.$$

Lemma. $(\mathbb{Z}/n\mathbb{Z})^\times$ equals the collection of congruence classes whose representatives are relatively prime to n :

$$(\mathbb{Z}/n\mathbb{Z})^\times = \{[a] \in \mathbb{Z}/n\mathbb{Z} \mid (a, n) = 1\}.$$

Axiom of Choice and Its Equivalences

Definition 2.21. Let (P, \leq) be a partially ordered set. Suppose $A \subset P$.

- (i) $u \in P$ is an **upper bound** for A if $x \leq u$ for all $x \in A$.
- (ii) $m \in P$ is a **maximal element** of P if $x \in P$ and $m \leq x$ implies $m = x$.
- (iii) Similarly we define **lower bound** and **minimal element**.
- (iv) $C \subset P$ is called a **chain** if either $x \leq y$ or $y \leq x$ for all $x, y \in C$.

This terminology of partially ordered sets will often be applied to an arbitrary family of sets. When this is done, it should be understood that the family is being regarded as a partially ordered set under the relation \subsetneq . Thus a maximal member of \mathcal{A} is a set $M \in \mathcal{A}$ such that M is a proper subset of no other member of \mathcal{A} ; a chain of sets is a family \mathcal{C} of sets such that $A \subsetneq B$ or $B \subsetneq A$ for all $A, B \in \mathcal{C}$.

Definition 2.22. Let \mathcal{F} be a family of sets. Then \mathcal{F} is said to be a *family of finite character* if for each set A , we have $A \in \mathcal{F}$ if and only if each finite subset of A is in \mathcal{F} .

We shall need the following technical fact.

Lemma 2.23. Let \mathcal{F} be a family of finite character, and let \mathcal{C} be a chain in \mathcal{F} . Then $\bigcup \mathcal{C} \in \mathcal{F}$.

Proof. It suffices to show that each finite subset of $\bigcup \mathcal{C}$ is in \mathcal{F} . Let $F = \{x_1, \dots, x_n\} \subset \bigcup \mathcal{C}$. Then there exist sets $C_1, \dots, C_n \in \mathcal{C}$ such that $x_i \in C_i$ ($i = 1, \dots, n$). Since \mathcal{C} is a chain, there exists $i_0 \in \{1, \dots, n\}$ such that $C_i \subsetneq C_{i_0}$ for $i = 1, \dots, n$. Then $F \subset C_{i_0} \in \mathcal{F}$. But \mathcal{F} is of finite character, and so $F \in \mathcal{F}$. \square

Theorem 2.24. The following are equivalent:

- (i) Axiom of choice: *The Cartesian product of any non-empty collection of non-empty sets is non-empty.*
- (ii) Tukey's lemma: *Every non-empty family of finite character has a maximal member.*
- (iii) Hausdorff maximality principle: *Every non-empty partially ordered set contains a maximal chain.*
- (iv) Zorn's lemma: *Every non-empty partially ordered set in which every chain has an upper bound has a maximal element.*
- (v) Well-ordering principle: *Every non-empty set has a well-ordering.*

Proof. We direct the reader to Section 3 of [HS65] for the complete proof. \square

Remark. It is a non-trivial result that Zorn's lemma is independent of the usual (Zermelo–Fraenkel) axioms of set theory in the sense that if the axioms of set theory are consistent, then so are these axioms together with Zorn's lemma; and if the axioms of set theory are consistent, then so are these axioms together with the negation of Zorn's lemma.

§2.3 Functions

Definitions

Definition 2.25 (Function). A **function** $f: X \rightarrow Y$ is a mapping of every element of X to some element of Y ; X and Y are known as the *domain* and *codomain* of f respectively.

Remark. The definition requires that a unique element of the codomain is assigned for every element of the domain. For example, for a function $f: \mathbb{R} \rightarrow \mathbb{R}$, the assignment $f(x) = \frac{1}{x}$ is not sufficient as it fails at $x = 0$. Similarly, $f(x) = y$ where $y^2 = x$ fails because $f(x)$ is undefined for $x < 0$, and for $x > 0$ it does not return a unique value; in such cases, we say the function is *ill-defined*. We are interested in the opposite; functions that are *well-defined*.

If a function is defined on some larger domain than we care about, it may be helpful to restrict the domain:

Definition 2.26 (Restriction). Suppose $f: X \rightarrow Y$. The **restriction** of f to $A \subset X$ is the map $f|_A: A \rightarrow Y$.

Remark. The restriction is almost the same function as the original function—just the domain has changed.

Another rather trivial but nevertheless important function is the identity map:

Definition 2.27 (Identity map). Given a set X , the **identity** $\text{id}_X: X \rightarrow X$ is defined by

$$\text{id}_X(x) = x \quad (x \in X).$$

Notation. If the domain is unambiguous, the subscript may be omitted.

Injectivity, Surjectivity, Bijectivity

Definition 2.28. Suppose $f: X \rightarrow Y$.

(i) f is **injective** (or *one-to-one*) if each element of Y has at most one element of X that maps to it:

$$\forall x_1, x_2 \in X, \quad f(x_1) = f(x_2) \implies x_1 = x_2$$

(ii) f is **surjective** (or *onto*) if every element of Y is mapped to at least one element of X :

$$\forall y \in Y, \quad \exists x \in X, \quad f(x) = y$$

(iii) f is **bijective** if it is both injective and surjective; a bijective function is termed a *bijection*.

Notation. We write $X \sim Y$ if there exists a bijection $f: X \rightarrow Y$.

Images and Pre-images

Definition 2.29. Suppose $f: X \rightarrow Y$. The **image** of $A \subset X$ under f is

$$f(A) := \{y \in Y \mid \exists x \in A, y = f(x)\}.$$

The **pre-image** of $B \subset Y$ under f is

$$f^{-1}(B) := \{x \in X \mid f(x) \in B\}.$$

Remark. Note the distinction between “codomain” and “range”.

Lemma 2.30. Let $f: X \rightarrow Y$. Suppose $A \subset X$ and $B \subset Y$.

- (i) If $A = f^{-1}(B)$, then $f(A) \subset B$.
- (ii) If $B = f(A)$, then $A \subset f^{-1}(B)$.

Proof.

- (i) Let $y = f(A)$. Then $y = f(x)$ for some $x \in A$. Since $A = f^{-1}(B)$, then $x \in f^{-1}(B)$. Then $f(x) = w$ for some $w \in B$. Thus $y = f(x) = w \in B$. Hence $f(A) \subset B$.
- (ii) Let $x \in A$. Then $f(x) \in f(A) = B$; let $f(x) = y$ for some $y \in B$. Consider $y \in B$; it could have one or more elements of A mapped to it. Hence $A \subset f^{-1}(B)$.

□

Remark. In general, we cannot conclude that $B = f(A)$ implies $A = f^{-1}(B)$.

We can express the previous result as follows:

$$f\left(f^{-1}(B)\right) \subset B, \quad A \subset f^{-1}\left(f(A)\right).$$

Lemma 2.31 (Algebra of pre-images). Suppose $f: X \rightarrow Y$. Then

- (i) $f^{-1}(A^c) = [f^{-1}(A)]^c$ for every $A \subset Y$;
- (ii) $f^{-1}\left(\bigcup_{i \in I} A_i\right) = \bigcup_{i \in I} f^{-1}(A_i)$;
- (iii) $f^{-1}\left(\bigcap_{i \in I} A_i\right) = \bigcap_{i \in I} f^{-1}(A_i)$.

Proof.

- (i) Suppose $A \subset Y$. Let $x \in X$, then

$$\begin{aligned} x \in f^{-1}(A^c) &\iff f(x) \in A^c \\ &\iff f(x) \notin A \\ &\iff x \notin f^{-1}(A) \\ &\iff x \in f^{-1}(A)^c \end{aligned}$$

Hence $f^{-1}(A^c) = f^{-1}(A)^c$.

(ii) Suppose $\{A_i \mid i \in I\}$ is a collection of subsets of Y . Then

$$\begin{aligned} x \in f^{-1}\left(\bigcup_{i \in I} A_i\right) &\iff f(x) \in \bigcup_{i \in I} A_i \\ &\iff f(x) \in A_i \text{ for some } i \in I \\ &\iff x \in f^{-1}(A_i) \text{ for some } i \in I \\ &\iff x \in \bigcup_{i \in I} f^{-1}(A_i) \end{aligned}$$

Hence $f^{-1}\left(\bigcup_{i \in I} A_i\right) = \bigcup_{i \in I} f^{-1}(A_i)$.

(iii) Suppose $\{A_i \mid i \in I\}$ is a collection of subsets of Y . Then

$$\begin{aligned} x \in f^{-1}\left(\bigcap_{i \in I} A_i\right) &\iff f(x) \in \bigcap_{i \in I} A_i \\ &\iff f(x) \in A_i \text{ for every } i \in I \\ &\iff x \in f^{-1}(A_i) \text{ for every } i \in I \\ &\iff x \in \bigcap_{i \in I} f^{-1}(A_i) \end{aligned}$$

Hence $f^{-1}\left(\bigcap_{i \in I} A_i\right) = \bigcap_{i \in I} f^{-1}(A_i)$.

□

Lemma 2.32 (Algebra of images). *Suppose $f: X \rightarrow Y$. Then*

- (i) $f(A)^c \subset f(A^c)$;
- (ii) $f\left(\bigcup_{i \in I} A_i\right) = \bigcup_{i \in I} f(A_i)$;
- (iii) $f\left(\bigcap_{i \in I} A_i\right) \subset \bigcap_{i \in I} f(A_i)$.

Composition

Definition 2.33 (Composition). Given $f: X \rightarrow Y$ and $g: Y \rightarrow Z$, the **composition** $g \circ f: X \rightarrow Z$ is defined by

$$(g \circ f)(x) = g(f(x)) \quad (\forall x \in X)$$

The composition of functions is not commutative. However, composition is associative, as the following results shows:

Proposition 2.34 (Associativity of composition). *Suppose $f: X \rightarrow Y$, $g: Y \rightarrow Z$, $h: Z \rightarrow W$. Then*

$$f \circ (g \circ h) = (f \circ g) \circ h.$$

Proof. Let $x \in X$. By the definition of composition, we have

$$(f \circ (g \circ h))(x) = f((g \circ h)(x)) = f(g(h(x))) = (f \circ g)(h(x)) = ((f \circ g) \circ h)(x).$$

□

Proposition 2.35 (Composition preserves injectivity and surjectivity).

- (i) If $f: X \rightarrow Y$ is injective and $g: Y \rightarrow Z$ is injective, then $g \circ f: X \rightarrow Z$ is injective.
(ii) If $f: X \rightarrow Y$ is surjective and $g: Y \rightarrow Z$ is surjective, then $g \circ f: X \rightarrow Z$ is surjective.

Proof.

- (i) Let $f: X \rightarrow Y$ and $g: Y \rightarrow Z$ be injective. To prove that $g \circ f: X \rightarrow Z$ is injective, we need to prove: for all $x, x' \in X$,

$$(g \circ f)(x) = (g \circ f)(x') \implies x = x'.$$

Suppose that $(g \circ f)(x) = (g \circ f)(x')$. Then by definition

$$g(f(x)) = g(f(x')).$$

Injectivity of g implies

$$f(x) = f(x'),$$

and injectivity of f implies

$$x = x'.$$

- (ii) Let $f: X \rightarrow Y$ and $g: Y \rightarrow Z$ be surjective. To prove that $g \circ f: X \rightarrow Z$ is surjective, we need to prove that for any $z \in Z$, there exists $x \in X$ such that $(g \circ f)(x) = z$.

Let $z \in Z$. By surjectivity of $g: Y \rightarrow Z$, there exists $y \in Y$ such that $g(y) = z$. By surjectivity of $f: X \rightarrow Y$, there exists $x \in X$ such that $f(x) = y$. This means that there exists $x \in X$ such that $g(f(x)) = g(y) = z$, as desired.

□

Proposition 2.36. $f: X \rightarrow Y$ is injective if and only if for any set Z and any functions $g_1, g_2: Z \rightarrow X$,

$$f \circ g_1 = f \circ g_2 \implies g_1 = g_2.$$

Proof.

\implies Suppose f is injective, and suppose $f \circ g_1 = f \circ g_2$. Let $z \in Z$. Then we have

$$f(g_1(z)) = f(g_2(z)).$$

Injectivity of f implies

$$g_1(z) = g_2(z),$$

so $g_1 = g_2$ (since the choice of $z \in Z$ is arbitrary).

\Leftarrow Pick $Z = \{1\}$, basically some random one-element set. Then for $x, y \in X$, define

$$\begin{aligned} g_1: Z &\rightarrow X, & g_1(1) &= x, \\ g_2: Z &\rightarrow Y, & g_2(1) &= y. \end{aligned}$$

Then for $x, y \in X$,

$$f(x) = f(y) \implies f(g_1(1)) = f(g_2(1)) \implies g_1(1) = g_2(1) \implies x = y$$

which shows that f is injective. \square

Proposition 2.37. $f: X \rightarrow Y$ is surjective if and only if for any set Z and any functions $g_1, g_2: Y \rightarrow Z$,

$$g_1 \circ f = g_2 \circ f \implies g_1 = g_2.$$

Proof.

\Rightarrow Suppose that f is surjective. Let $y \in Y$. Surjectivity of f means there exists $x \in X$ such that $f(x) = y$. Then

$$g_1 \circ f = g_2 \circ f \implies g_1(f(x)) = g_2(f(x)) \implies g_1(y) = g_2(y)$$

so $g_1 = g_2$.

\Leftarrow We prove the contrapositive. Suppose f is not surjective, then there exists $y \in Y$ such that for all $x \in X$ we have $f(x) \neq y$. We then aim to construct set Z and $g_1, g_2: Y \rightarrow Z$ such that

$$(i) \quad g_1(y) \neq g_2(y)$$

$$(ii) \quad \forall y' \neq y, g_1(y') = g_2(y')$$

Because if this is satisfied, then $\forall x \in X$, since $f(x) \neq y$ we have from (ii) that $g_1(f(x)) = g_2(f(x))$; thus $g_1 \circ f = g_2 \circ f$, and yet from (i) we have $g_1 \neq g_2$.

We construct $Z = Y \cup \{1, 2\}$ for some random $1, 2 \notin Y$.

Then we define

$$\begin{aligned} g_1: Y &\rightarrow Z, & g_1(y) &= 1, & g_1(y') &= y' \\ g_2: Y &\rightarrow Z, & g_2(y) &= 2, & g_2(y') &= y' \end{aligned}$$

Then when y is not in the image of f , these two functions will satisfy $g_1 \circ f = g_2 \circ f$ but not $g_1 = g_2$.

So conversely, if for any set Z and any functions $g_i: Y \rightarrow Z$ we have $g_1 \circ f = g_2 \circ f \implies g_1 = g_2$, such a value y that is in the codomain but not in the range of f cannot appear, and hence f must be surjective. \square

Lemma 2.38 (Inverse image of composition). Suppose $f: X \rightarrow Y$, $g: Y \rightarrow Z$. Then

$$(g \circ f)^{-1}(A) = f^{-1}(g^{-1}(A))$$

for every $A \subset Z$.

Proof. Suppose $A \subset Z$. Let $x \in X$, then we have

$$\begin{aligned} x \in (g \circ f)^{-1}(A) &\iff (g \circ f)(x) \in A \\ &\iff g(f(x)) \in A \\ &\iff f(x) \in g^{-1}(A) \\ &\iff x \in f^{-1}(g^{-1}(A)) \end{aligned}$$

Hence $(g \circ f)^{-1}(A) = f^{-1}(g^{-1}(A))$. □

Invertibility

Recalling that id_X is the identity map on X , we can define invertibility.

Definition 2.39 (Invertibility). Suppose $f: X \rightarrow Y$. We say that

- (i) f is **left-invertible** if there exists $g: Y \rightarrow X$ such that $g \circ f = \text{id}_X$; we call g a *left-inverse* of f ;
- (ii) f is **right-invertible** if there exists $h: Y \rightarrow X$ such that $f \circ h = \text{id}_Y$; we call h a *right-inverse* of f ;
- (iii) f is **invertible** if there exists $k: Y \rightarrow X$ which is a left and right inverse of f ; we call k an *inverse* of f .

Remark. Notice that if g is left-inverse to f then f is right-inverse to g . A function can have more than one left-inverse, or more than one right-inverse.

Example 2.40. Let

$$\begin{aligned} f: \mathbb{R} &\rightarrow [0, \infty), & f(x) &= x^2 \\ g: [0, \infty) &\rightarrow \mathbb{R}, & g(x) &= \sqrt{x} \end{aligned}$$

- f is not left-invertible. Suppose otherwise, for a contradiction, that h is a left inverse of f , so that $hf = \text{id}_{\mathbb{R}}$. Then

Lemma 2.41 (Uniqueness of inverse). *If $f: X \rightarrow Y$ is invertible, then its inverse is unique.*

Proof. Let g_1 and g_2 be two functions for which $g_i \circ f = \text{id}_X$ and $f \circ g_i = \text{id}_Y$. Using the fact that composition is associative, and the definition of the identity maps, we can write

$$g_1 = g_1 \circ \text{id}_Y = g_1 \circ (f \circ g_2) = (g_1 \circ f) \circ g_2 = \text{id}_X \circ g_2 = g_2.$$

□

Since the inverse is unique, we can give it a notation.

Notation. The inverse of f is denoted by f^{-1}

Remark. Immediately from the definition, if f is invertible then f^{-1} is also invertible, and $(f^{-1})^{-1} = f$.

The following result provides an important and useful criterion for invertibility.

Lemma 2.42 (Invertibility criterion). *Suppose $f: X \rightarrow Y$. Then*

- (i) *f is left-invertible if and only if f is injective;*
- (ii) *f is right-invertible if and only if f is surjective;*
- (iii) *f is invertible if and only if f is bijective.*

Proof.

(i) \Rightarrow Suppose f is left-invertible; let g be a left-inverse of f , so $g \circ f = \text{id}_X$.

Now suppose $f(a) = f(b)$. Then applying g to both sides gives $g(f(a)) = g(f(b))$, so $a = b$.

\Leftarrow Let f be injective. Choose any x_0 in the domain of f . Define $g: Y \rightarrow X$ as follows; note that each $y \in Y$ is either in the image of f or not.

- If y is in the image of f , it equals $f(x)$ for a *unique* $x \in X$ (uniqueness is because of the injectivity of f), so define $g(y) = x$.
- If y is not in the image of f , define $g(y) = x_0$.

Clearly $g \circ f = \text{id}_X$.

(ii) \Rightarrow Suppose f is right-invertible; let g be a right-inverse of f , so $f \circ g = \text{id}_Y$.

Let $y \in Y$. Then $f(g(y)) = \text{id}_Y(y) = y$ so $y \in f(X)$. Thus $f(X) = Y$ so f is surjective.

\Leftarrow Suppose f is surjective. Let $y \in Y$, then y is in the image of f , so we can choose an element $g(y) \in X$ such that $f(g(y)) = y$. This defines a function $g: Y \rightarrow X$ which is evidently a right-inverse of f .

(iii) \Rightarrow Suppose f is invertible. Then f is left-invertible and right-invertible. By (i) and (ii), f is injective and surjective, so f is bijective.

\Leftarrow Suppose f is bijective. Then by (i) and (ii), f has a left-inverse $g: Y \rightarrow X$ and a right-inverse $h: Y \rightarrow X$. But “invertible” requires a single function to be *both* a left and right inverse, so we need to show that $g = h$:

$$g = g \circ \text{id}_Y = g \circ (f \circ h) = (g \circ f) \circ h = \text{id}_X \circ h = h$$

so $g = h$ is an inverse of f .

□

The following result shows how to invert the composition of invertible functions.

Proposition 2.43 (Inverse of composition). *Suppose $f: X \rightarrow Y$, $g: Y \rightarrow Z$. If f and g are invertible, then $g \circ f$ is invertible, and*

$$(g \circ f)^{-1} = f^{-1} \circ g^{-1}.$$

Proof. Making repeated use of the fact that function composition is associative, and the definition of the inverses f^{-1} and g^{-1} , we note that

$$\begin{aligned}(f^{-1} \circ g^{-1}) \circ (g \circ f) &= ((f^{-1} \circ g^{-1}) \circ g) \circ f \\ &= (f^{-1} \circ (g^{-1} \circ g)) \circ f \\ &= (f^{-1} \circ \text{id}_Y) \circ f \\ &= f^{-1} \circ f \\ &= \text{id}_X\end{aligned}$$

and similarly,

$$\begin{aligned}(g \circ f) \circ (f^{-1} \circ g^{-1}) &= g \circ (f \circ (f^{-1} \circ g^{-1})) \\ &= g \circ ((f \circ f^{-1}) \circ g^{-1}) \\ &= g \circ (\text{id}_X \circ g^{-1}) \\ &= g \circ g^{-1} \\ &= \text{id}_Z\end{aligned}$$

which shows that $f^{-1} \circ g^{-1}$ satisfies the properties required to be the inverse of $g \circ f$. \square

Corollary 2.44. *If f_1, \dots, f_n are invertible and the composition $f_1 \circ \dots \circ f_n$ makes sense, then it is also invertible and its inverse is*

$$f_n^{-1} \circ \dots \circ f_1^{-1}.$$

Proposition 2.45. *\sim is an equivalence relation between sets.*

Proof. We need to prove (i) reflexivity, (ii) symmetry, and (iii) transitivity.

- (i) The identity map gives a bijection from a set to itself.
- (ii) Suppose $f: X \rightarrow Y$ is a bijection. Then f is invertible, with inverse $f^{-1}: Y \rightarrow X$. Since f^{-1} is invertible (with inverse f), it is bijective.
- (iii) Suppose $f: X \rightarrow Y$ and $g: Y \rightarrow Z$ are bijections, and thus they are invertible. Then by the previous result, $g \circ f$ is invertible and thus bijective.

\square

Theorem 2.46 (Cantor–Schröder–Bernstein). *If $f: X \rightarrow Y$ and $g: Y \rightarrow X$ are injective, then $A \sim B$.*

§2.4 Cardinality

This section is about formalising the notion of the “size” of a set.

Definition 2.47. A and B said to be *equivalent* (or have the same *cardinality*), denoted by $A \sim B$, if there exists a bijection $f: A \rightarrow B$.

Notation. For $n \in \mathbb{N}$, denote

$$\begin{aligned}\mathbb{N}_n &= \{k \in \mathbb{N} \mid 1 \leq k \leq n\}, \\ n\mathbb{N} &= \{nk \mid k \in \mathbb{N}\}.\end{aligned}$$

Definition 2.48. For any set A , we say

- (i) A is *finite* if $A \sim \mathbb{N}_n$ for some integer $n \in \mathbb{N}$, then the *cardinality* of A is $|A| = n$; A is *infinite* if A is not finite;
- (ii) A is *countable* if $A \sim \mathbb{N}$; A is *uncountable* if A is neither finite nor countable; A is *at most countable* if A is finite or countable.

Remark. Any countable set can be “listed” in a sequence a_1, a_2, \dots of distinct terms. This technique is particularly useful when there is not possible to deduce an explicit formula for a bijection.

Lemma 2.49. \mathbb{N} is infinite.

Proof. We want to show that there does not exist a bijection from \mathbb{N}_n to \mathbb{N} , for all $n \in \mathbb{N}$. We prove by induction on n .

For the base case $n = 1$, if there exists a function $f_1: \{1\} \rightarrow \mathbb{N}$, consider the set $\mathbb{N} \setminus f_1(\{1\})$. It is not empty, so f_1 is not surjective, thus it is not bijective.

For the inductive step, we want to show if there does not exist a bijection from \mathbb{N}_k to \mathbb{N} , then there does not exist a bijection from \mathbb{N}_{k+1} to \mathbb{N} . We prove the contrapositive: if there exists a bijection from $\mathbb{N}_{k+1} \rightarrow \mathbb{N}$, then there exists a bijection from \mathbb{N}_k to \mathbb{N} .

Suppose $h: \mathbb{N}_{k+1} \rightarrow \mathbb{N}$ is a bijection. If remove the element $k + 1$, then there exists a bijection from \mathbb{N}_k to $\mathbb{N} \setminus \{h(k + 1)\}$. But $\mathbb{N} \setminus \{h(k + 1)\} \sim \mathbb{N}$ so $\mathbb{N}_k \sim \mathbb{N}$. □

Corollary 2.50. Any countable set is infinite.

Example 2.51. \mathbb{N} is countable since the identity map from \mathbb{N} to \mathbb{N} is a bijection.

Example 2.52. $n\mathbb{N}$ is countable.

Proof. Let $f: \mathbb{N} \rightarrow n\mathbb{N}$ which sends $k \mapsto nk$. We need to show that f is bijective:

- For any $k_1, k_2 \in \mathbb{N}$, $nk_1 = nk_2$ implies $k_1 = k_2$ so f is injective.
- For any $x \in n\mathbb{N}$, $x = nk$ for some $k \in \mathbb{N}$, thus $\frac{x}{n} = k \in \mathbb{N}$ so f is surjective.

Hence f is bijective, so $n\mathbb{N} \sim \mathbb{N}$ and we are done. □

Example 2.53. \mathbb{Z} is countable.

Proof. Consider the following arrangement of the elements of \mathbb{Z} and \mathbb{N} :

$$\mathbb{Z}: \quad 0, 1, -1, 2, -2, 3, -3, \dots$$

$$\mathbb{N}: \quad 1, 2, 3, 4, 5, 6, 7, \dots$$

In fact we can write an explicit formula for a bijection $f: \mathbb{N} \rightarrow \mathbb{Z}$ where

$$f(n) = \begin{cases} \frac{n}{2} & (n \text{ even}) \\ -\frac{n-1}{2} & (n \text{ odd}) \end{cases}$$

□

Proposition 2.54. *Every infinite subset of a countable set is countable.*

Proof. Let S be the countable set. Then we can arrange the elements of S in a sequence (s_n) of distinct elements:

$$s_1, s_2, \dots$$

Suppose $E \subset S$ is infinite. The main idea is to show that we can list out the elements of E in a sequence. We now construct a sequence (n_k) as follows: Let

$$n_1 = \min\{i \mid s_i \in E\}$$

$$n_2 = \min\{i \mid s_i \in E, i > n_1\}$$

$$\vdots$$

$$n_k = \min\{i \mid s_i \in E, i > n_{k-1}\}.$$

Then

$$E = \{s_{n_1}, s_{n_2}, \dots\},$$

where we note that the function $f(k) = s_{n_k}$ ($k = 1, 2, \dots$) is bijective. Hence $E \sim \mathbb{N}$, as desired. □

Remark. This shows that countable sets represent the “smallest” infinity: No uncountable set can be a subset of a countable set.

Proposition 2.55. *The countable union of countable sets is countable.*

Proof. Let $\{A_n \mid n \in \mathbb{N}\}$ be a family of countable sets; clearly this is a countable collection of sets (indexed by \mathbb{N}). Then we want to show that the union

$$S = \bigcup_{n=1}^{\infty} A_n$$

is countable.

Since every set A_n is countable, we can list its elements in a sequence (a_{nk}) ($k = 1, 2, 3, \dots$). Arrange the elements of all the sets in $\{A_n\}$ in the form of an infinite array, containing all elements of S , where the elements of A_n form the n -th row.

$$\begin{array}{l}
 A_1: \quad \cancel{a_{11}} \quad \cancel{a_{12}} \quad \cancel{a_{13}} \quad \cancel{a_{14}} \quad \cdots \\
 A_2: \quad \cancel{a_{21}} \quad \cancel{a_{22}} \quad \cancel{a_{23}} \quad \cancel{a_{24}} \quad \cdots \\
 A_3: \quad \cancel{a_{31}} \quad \cancel{a_{32}} \quad \cancel{a_{33}} \quad \cancel{a_{34}} \quad \cdots \\
 A_4: \quad \cancel{a_{41}} \quad \cancel{a_{42}} \quad \cancel{a_{43}} \quad \cancel{a_{44}} \quad \cdots \\
 \vdots
 \end{array}$$

We then zigzag our way through the array, and arrange these elements in a sequence

$$a_{11}, a_{21}, a_{12}, a_{31}, a_{22}, a_{13}, a_{41}, a_{32}, a_{23}, a_{14}, \dots$$

thus S is countable, and we are almost done!

A small problem is that if any two of the sets A_n have elements in common, these will appear more than once in the above sequence. Then we take a subset $T \subset S$, where every element only appears once. Note that T is an infinite subset, since $A_1 \subset T$ is infinite. Then since T is an infinite subset of a countable set S , by Proposition 2.54, T is countable. \square

Remark. If we were to instead start by going down by the first row of the above array, then we would not get to the second row (and beyond); all that would show is the first row is countable. Instead, we wind our way through diagonally, ensuring that we hit every number of the array.

Corollary 2.56. *Suppose A is an indexing set that is at most countable. Let $\{B_\alpha \mid \alpha \in A\}$ be a family of sets that are at most countable. Then the union*

$$\bigcup_{\alpha \in A} B_\alpha$$

is at most countable.

Proposition 2.57. *Let A be a countable set. For $n \in \mathbb{N}$, let*

$$B_n = \{(a_1, \dots, a_n) \mid a_i \in A\}.$$

Then B_n is countable.

Proof. We prove by induction on n . That B_1 is countable is evident, since $B_1 = A$.

Now suppose B_{n-1} is countable. The elements of B_n are of the form

$$(b, a) \quad (b \in B_{n-1}, a \in A)$$

For every fixed b , the set of ordered pairs (b, a) is equivalent to A , and hence countable. Thus B_n is a union of countable sets. By Proposition 2.55, B_n is countable. \square

Corollary 2.58. \mathbb{Q} is countable.

Proof. Note that every $x \in \mathbb{Q}$ is of the form $\frac{b}{a}$, where $a, b \in \mathbb{Z}$. By the previous result, taking $n = 2$, the set of pairs (a, b) and therefore the set of fractions $\frac{b}{a}$ is countable. \square

That not all infinite sets are, however, countable, is shown by the next result.

Proposition 2.59. *Let A be the set of all sequences whose elements are the digits 0 and 1. Then A is uncountable.*

Proof. Let $E \subset A$ be countable, consisting of the sequences s_1, s_2, s_3, \dots .

We construct a new sequence s as follows:

$$n\text{-th digit of } s = \begin{cases} 0 & \text{if } n\text{-th digit in } s_n \text{ is 1,} \\ 1 & \text{if } n\text{-th digit in } s_n \text{ is 0.} \end{cases}$$

Then the sequence s differs from every member of E in at least one place, so $s \notin E$. But clearly $s \in A$; hence $E \subsetneq A$.

We have shown that every countable subset of A is a proper subset of A . It follows that A is uncountable (for otherwise A would be a proper subset of A , which is absurd). \square

Remark. The idea of the above proof is called *Cantor's diagonal process*, first used by Cantor. This is because if elements of the sequences s_1, s_2, s_3, \dots are listed out in an array, it is the elements on the diagonal which are involved in the construction of the new sequence.

Corollary 2.60. \mathbb{R} is uncountable.

Proof. This follows from the binary representation of the real numbers. \square

Theorem 2.61 (Cantor's theorem). *For any set A , we have $A \not\sim \mathcal{P}(A)$.*

Proof. Suppose otherwise, for a contradiction, that $A \sim \mathcal{P}(A)$. Then there exists a bijection $f: A \rightarrow \mathcal{P}(A)$. Then for each $x \in A$, $f(x)$ is a subset of A . Now consider the "anti-diagonal" set

$$B = \{x \in A \mid x \notin f(x)\}.$$

That is, B is the subset of A containing all $x \in A$ such that x is not in the set $f(x)$. Since $B \subset A$, we have $B \in \mathcal{P}(A)$. Since f is bijective (in particular surjective), there exists $x \in A$ such that $f(x) = B$. Now there are two cases: (i) $x \in B$, or (ii) $x \notin B$.

(i) If $x \in B$, then by definition of the set B it must be the case that $x \notin f(x)$. But since $f(x) = B$, we then have $x \notin B$. This is absurd since we cannot have $x \in B$ and $x \notin B$ simultaneously.

(ii) If $x \notin B$, by definition of the set B , this implies that $x \in f(x)$. But $f(x) = B$. So we have $x \in B$ and $x \notin B$, which is again absurd.

In either case, we have reached a contradiction. Hence there cannot exist a surjective (and thus bijective) function $A \rightarrow \mathcal{P}(A)$. \square

Exercises

Exercise 2.1. Prove that the statement $\forall x \in \emptyset, P(x)$ is vacuously true.

Solution. Let S be the embedding set. The statement $\forall x \in \emptyset, P(x)$ means

$$\forall x \in S, \quad x \in \emptyset \implies P(x).$$

But $x \in \emptyset$ is always false, by the definition of empty set. Hence the statement is always true, regardless of x . \square

Exercise 2.2. Prove that for any set $A \subset S$, $\emptyset \subset A$ and $A \subset A$.

Solution. Let $A \subset S$. Let $x \in \emptyset$, then $x \in \emptyset \implies x \in A$ is vacuously true, so $\emptyset \subset A$.

Likewise, let $x \in A$, then $x \in A \implies x \in A$ is always true, so $A \subset A$. \square

Exercise 2.3. Let A be the set of all complex polynomials in n variables. Given a subset $T \subset A$, define the *zeros* of T as the set

$$Z(T) = \{P = (a_1, \dots, a_n) \mid f(P) = 0 \text{ for all } f \in T\}.$$

$Y \subset \mathbb{C}^n$ is called an *algebraic set* if there exists $T \subset A$ such that $Y = Z(T)$.

Prove that the union of two algebraic sets is an algebraic set.

Solution. We would like to consider $T = \{f_1, f_2, \dots\}$ expressed as indexed sets $T = \{f_i\}$. Then $Z(T)$ can also be expressed as $\{P \mid \forall i, f_i(P) = 0\}$.

Suppose that we have two algebraic sets X and Y . Let $X = Z(S)$, $Y = Z(T)$ where S, T are subsets of A (basically, they are certain sets of polynomials). Then

$$X = \{P \mid \forall f \in S, f(P) = 0\}$$

$$Y = \{P \mid \forall g \in T, g(P) = 0\}$$

We imagine that for $P \in X \cap Y$, we have $f(P) = 0$ or $g(P) = 0$. Hence we consider the set of polynomials

$$U = \{f \cdot g \mid f \in S, g \in T\}$$

For any $P \in X \cup Y$ and for any $fg \in U$ where $f \in S$ and $f \in g$, either $f(P) = 0$ or $g(P) = 0$, hence $fg(P) = 0$ and thus $P \in Z(U)$.

On the other hand if $P \in Z(U)$, suppose otherwise that P is not in $X \cup Y$, then P is neither in X nor in Y . This means that there exists $f \in S, g \in T$ such that $f(P) \neq 0$ and $g(P) \neq 0$, hence $fg(P) \neq 0$. This is a contradiction as $P \in Z(U)$ implies $fg(P) = 0$. Hence we have $X \cup Y = Z(U)$ and thus $X \cup Y$ is an algebraic set.

Now the other direction is simpler and can actually be generalised: The intersection of arbitrarily many algebraic sets is algebraic.

The basic result is that if $X = Z(S)$ and $Y = Z(T)$ then $X \cap Y = Z(S \cup T)$. \square

Exercise 2.4. Let $A = \mathbb{R}$ and for any $x, y \in A$, $x \sim y$ if and only if $x - y \in \mathbb{Z}$. For any two equivalence classes $[x], [y] \in A/\sim$, define

$$[x] + [y] = [x + y] \text{ and } -[x] = [-x]$$

- (a) Show that the above definitions are well-defined.
 (b) Find a one-to-one correspondence $\phi: X \rightarrow Y$ between $X = A/\sim$ and $Y: |z| = 1$, i.e. the unit circle in \mathbb{C} , such that for any $[x_1], [x_2] \in X$ we have

$$\phi([x_1])\phi([x_2]) = \phi([x_1 + x_2])$$

- (c) Show that for any $[x] \in X$,

$$\phi(-[x]) = \phi([x])^{-1}$$

Solution.

- (a)

$$(x' + y') - (x + y) = (x' - x) + (y' - y) \in \mathbb{Z}$$

Thus $[x' + y'] = [x + y]$

$$(-x') - (-x) = -(x' - x) \in \mathbb{Z}$$

Thus $[-x'] = [-x]$.

- (b) Complex numbers in the polar form: $z = re^{i\theta}$

Then the correspondence is given by $\phi([x]) = e^{2\pi ix}$

$$[x] = [y] \iff x - y \in \mathbb{Z} \iff e^{2\pi i(x-y)} = 1 \iff e^{2\pi ix} = e^{2\pi iy}$$

Hence this is a bijection.

Before that, we also need to show that ϕ is well-defined, which is almost the same as the above.

If we choose another representative x' then

$$\phi([x]) = e^{2\pi ix'} = e^{2\pi ix} \cdot e^{2\pi i(x'-x)} = e^{2\pi ix}$$

- (c) You can either refer to the specific correspondence $\phi([x]) = e^{2\pi ix}$ or use its properties.

$$\phi(-[x])\phi([x]) = \phi([-x])\phi([x]) = \phi([-x + x]) = \phi([0]) = 1$$

□

Exercise 2.5 (Complex Numbers). Let $\mathbb{R}[x]$ denote the set of real polynomials. Define

$$\mathbb{C} = \mathbb{R}[x]/(x^2 + 1)\mathbb{R}[x]$$

where

$$f(x) \sim g(x) \iff x^2 + 1 \text{ divides } f(x) - g(x).$$

The complex number $a + bi$ is defined to be the equivalence class of $a + bx$.

- (a) Define the sum and product of two complex numbers and show that such definitions are well-defined.
- (b) Define the reciprocal of a complex number.

Exercise 2.6 ([Rud76] 2.2). $z \in \mathbb{C}$ is said to be *algebraic* if there exist integers a_0, \dots, a_n , not all zero, such that

$$a_0 z^n + a_1 z^{n-1} + \dots + a_{n-1} z + a_n = 0.$$

Prove that the set of all algebraic numbers is countable. *Hint*: For every positive integer N there are only finitely many equations with

$$n + |a_0| + |a_1| + \dots + |a_n| = N.$$

Solution. Following the hint, let A_N be the set of numbers z that satisfy $a_0 z^n + a_1 z^{n-1} + \dots + a_{n-1} z + a_n = 0$, for some coefficients a_0, \dots, a_n which satisfy

$$n + |a_0| + |a_1| + \dots + |a_n| = N.$$

By the fundamental theorem of algebra, $a_0 z^n + a_1 z^{n-1} + \dots + a_{n-1} z + a_n = 0$ has at most n solutions, so each A_N is finite. Hence the set of algebraic numbers, which is the union

$$\bigcup_{N=2}^{\infty} A_N$$

is at most countable. Since all rational numbers are algebraic, it follows that the set of algebraic numbers is exactly countable. \square

Exercise 2.7 ([Rud76] 2.3). Prove that there exist real numbers which are not algebraic.

Solution. By the previous exercise, the set of real algebraic numbers is countable. If every real number were algebraic, the entire set of real numbers would be countable, a contradiction. \square

Exercise 2.8 ([Rud76] 2.4). Is the set of irrational real numbers countable?

Solution. No. If $\mathbb{R} \setminus \mathbb{Q}$ were countable, $\mathbb{R} = \mathbb{Q} \cup (\mathbb{R} \setminus \mathbb{Q})$ would be countable, which is clearly false. \square

II

Abstract Algebra

Algebra is the study of collections of objects (sets, groups, rings, fields, etc). In algebra, we are concerned about the structures of these collections and how these collections interact than about the objects themselves. In fact, with homomorphism and isomorphisms, the original objects become irrelevant.

3 Groups

Summary

- Group, Abelian group, examples. Subgroups. subgroup generated by a subset of a group. Cyclic subgroups.
- Cosets and Lagrange's theorem; examples. The order of an element. Fermat's little theorem.
- Isomorphisms, examples. Groups of order 8 or less up to isomorphism (stated without proof). Homomorphisms of groups with motivating examples. Kernels. Images. Normal subgroups. Quotient groups; examples. First Isomorphism Theorem. Cayley's theorem.
- Group actions; examples. Definition of orbits and stabilizers. Transitivity. Orbits partition the set. Orbit-stabilizer Theorem. Examples and applications including Cauchy's Theorem and to conjugacy classes. Orbit-counting formula.

§3.1 Definition and Properties

Definition 3.1. A *binary operation* on a set G is a map $*$: $G \times G \rightarrow G$.

Notation. For any $a, b \in G$, if the operation is clear, we write ab for the image of (a, b) under $*$.

$*$ is *associative* if $(ab)c = a(bc)$ for all $a, b, c \in G$; $*$ is *commutative* if $ab = ba$ for all $a, b \in G$.

Definition 3.2 (Group). A *group* $(G, *)$ consists of a set G and a binary operation $*$ on G satisfying the following group axioms:

(i) $a(bc) = (ab)c$ for all $a, b, c \in G$; (associativity)

(ii) there exists $e \in G$ such that $ae = ea = a$ for all $a \in G$; (identity)

(iii) for all $a \in G$, there exists $c \in G$ such that $ac = ca = e$. (invertibility)

We say G is *abelian* if the operation is commutative; otherwise, G is *non-abelian*.

Remark. When verifying that $(G, *)$ is a group we have to check (i), (ii), (iii) above and also that $*$ is a binary operation closed in G —that is, $a * b \in G$ for all $a, b \in G$.

Notation. We simply denote a group $(G, *)$ by G if the operation is clear.

Notation. Since $*$ is associative, we omit unnecessary parentheses and write $(ab)c = a(bc) = abc$.

Notation. For any $a \in G$, $n \in \mathbb{Z}^+$, denote $a^n = \underbrace{a \cdot a \cdots a}_{n \text{ times}}$.

Notation. Denote the additive group $\mathbb{C} = (\mathbb{C}, +)$ etc., the multiplicative group $\mathbb{C}^\times = \mathbb{C} \setminus \{0\}$ etc., the set of (congruence classes of) integers modulo n under addition as $\mathbb{Z}/n\mathbb{Z}$ and under multiplication as $(\mathbb{Z}/n\mathbb{Z})^\times$.

Example 3.3.

- $\mathbb{Z}, \mathbb{Q}, \mathbb{R}, \mathbb{C}$ are groups, with identity 0 and (additive) inverse $-a$ for all a .

- $\mathbb{Q}^\times, \mathbb{R}^\times, \mathbb{C}^\times, \mathbb{Q}^+, \mathbb{R}^+$ are groups under \times , with identity 1 and (multiplicative) inverse $\frac{1}{a}$ for all a .
- For $n \in \mathbb{Z}^+$, $\mathbb{Z}/n\mathbb{Z}$ is an abelian group under $+$.
- For $n \in \mathbb{Z}^+$, $(\mathbb{Z}/n\mathbb{Z})^\times$ is an abelian group under multiplication.

Lemma 3.4. *Let G be a group.*

- (i) *The identity of G is unique.*
- (ii) *For each $a \in G$, a^{-1} is unique.*
- (iii) *$(a^{-1})^{-1} = a$ for all $a \in G$.*
- (iv) *$(ab)^{-1} = b^{-1}a^{-1}$.*
- (v) *For any $a_1, \dots, a_n \in G$, $a_1 \cdots a_n$ is independent of how we arrange the parantheses (generalised associative law).*

Proof.

- (i) Suppose that e and e' are identities of G . Then

$$e = ee' = e'$$

where the first equality holds as e' is an identity, and the second equality holds as e is an identity. Since $e = e'$, the identity is unique.

Notation. Denote the identity of G by 1_G ; the subscript may be omitted if there is no ambiguity.

- (ii) Suppose that b and c are both inverses of a . Then $ab = 1$, $ca = 1$, so

$$c = c1 = c(ab) = (ca)b = 1b = b.$$

Hence the inverse is unique.

Notation. Denote the inverse of $a \in G$ by a^{-1} .

- (iii) To show $(a^{-1})^{-1} = a$ is exactly the problem of showing that a is the inverse of a^{-1} , which is by definition of the inverse (with the roles of a and a^{-1} interchanged).
- (iv) Let $c = (ab)^{-1}$. Then $(ab)c = 1$, or $a(bc) = 1$ by associativity, which gives $bc = a^{-1}$. Applying b^{-1} on both sides gives $c = b^{-1}a^{-1}$.
- (v) The result is trivial for $n = 1, 2, 3$. For all $k < n$ assume that any $a_1 \cdots a_k$ is independent of parantheses. Then

$$(a_1 \cdots a_n) = (a_1 \cdots a_k)(a_{k+1} \cdots a_n).$$

By inductive hypothesis, both terms are independent of parentheses since $k, n - k < n$. Hence by induction we are done.

□

Lemma 3.5 (Cancellation law). *Let $a, b \in G$. Then the equations $ax = b$ and $ya = b$ have unique solutions for $x, y \in G$. In particular, we can cancel on the left and right.*

Proof. We can solve $ax = b$ by applying a^{-1} to both sides of the equation to get $x = a^{-1}b$. The uniqueness of x follows because a^{-1} is unique. A similar case holds for $ya = b$. \square

Definition 3.6 (Order of a group). Let G be a group. The **order** of G is its cardinality $|G|$. A group G is a *finite group* if $|G| < \infty$.

One way to represent a finite group is by means of a **Cayley table**. Let $G = \{1, g_2, g_3, \dots, g_n\}$. The Cayley table of G is a square grid which contains all the possible products of two elements from G : the product $g_i g_j$ appears in the i -th row and j -th column.

Remark. Note that a group is abelian if and only if its Cayley table is symmetric about the main (top-left to bottom-right) diagonal.

Examples

Example 3.7 (Dihedral groups). An important family of groups is the **dihedral groups**. For $n \in \mathbb{Z}^+$, $n \geq 3$, let D_{2n} be the set of symmetries of a regular n -gon.

Let r be the rotation clockwise about the origin by $\frac{2\pi}{n}$ radians, s be the reflection about the line of symmetry through the first labelled vertex and the origin. (Read from right to left: for instance, sr means do r then s .)

Properties of D_{2n} :

- $1, r, r^2, \dots, r^{n-1}$ are all distinct and $r^n = 1$, so $|r| = n$.
- $s^2 = 1$ since we either reflect or do not reflect, so $|s| = 2$.
- $s \neq r^i$ for any i , since the effect of any reflection cannot be obtained from any form of rotation.
- $sr^i \neq sr^j$ for all $i \neq j$ ($0 \leq i, j \leq n-1$), so

$$D_{2n} = \{1, r, \dots, r^{n-1}, s, sr, \dots, sr^{n-1}\}$$

and thus $|D_{2n}| = 2n$.

- $rs = sr^{-1}$
- $r^i s = sr^{-i}$

Proof: From above, this is true for $i = 1$. Assume it holds for $k < n$. Then $r^{k+1}s = r(r^k s) = r sr^{-k}$. Then $rs = sr^{-1}$ so $r sr^{-k} = sr^{-1} r^{-k} = sr^{-k-1}$ so we are done.

Note that for each $n \in \mathbb{Z}^+$, the generators of D_{2n} are r and s , and we have shown that they satisfy $r^n = 1$, $s^2 = 1$, and $rs = sr^{-1}$; these are called *relations*. Any other equation involving the generators can be derived from these relations.

Any such collection of generators S and relations R_1, \dots, R_m for a group G is called a *presentation*, written

$$G = \langle S \mid R_1, \dots, R_m \rangle.$$

For example,

$$D_{2n} = \langle r, s \mid r^n = s^2 = 1, rs = sr^{-1} \rangle.$$

Example 3.8 (Permutation groups). Let S be a non-empty set. A bijection $S \rightarrow S$ is called a *permutation* of S ; the set of permutations of S is denoted by $\text{Sym}(S)$.

We now show that $\text{Sym}(S)$ is a group under function composition \circ ; $(\text{Sym}(S), \circ)$ is the *symmetric group* on S . Note that \circ is a binary operation on $\text{Sym}(S)$ since if $\sigma: S \rightarrow S$ and $\tau: S \rightarrow S$ are both bijections, then $\sigma \circ \tau$ is also a bijection from S to S .

- (i) Function composition is associative so \circ is associative.
- (ii) The identity of $\text{Sym}(S)$ is the identity map 1 , defined by $1(a) = a$ for all $a \in S$.
- (iii) For every permutation σ , σ is bijective and thus invertible, so there exists a (2-sided) inverse $\sigma^{-1}: S \rightarrow S$ satisfying $\sigma \circ \sigma^{-1} = \sigma^{-1} \circ \sigma = 1$.

In the special case where $S = \{1, 2, \dots, n\}$, the symmetric group on S is called the *symmetric group of degree n* , denoted by S_n .

Proposition. *If $|S| \geq 3$ then $\text{Sym}(S)$ is non-abelian.*

Proof. Let $S = \{x_1, x_2, x_3\}$ where three elements are distinct. □

Proposition. $|S_n| = n!$

Proof. Obvious, since there are $n!$ permutations of $\{1, 2, \dots, n\}$. □

Example 3.9 (Matrix groups). For $n \in \mathbb{Z}^+$, let $GL_n(\mathbf{F})$ be the set of all $n \times n$ invertible matrices whose entries are in \mathbf{F} :

$$GL_n(\mathbf{F}) = \{A \in M_{n \times n}(\mathbf{F}) \mid \det(A) \neq 0\}.$$

We show that $GL_n(\mathbf{F})$ is a group under matrix multiplication; $GL_n(\mathbf{F})$ is the *general linear group* of degree n .

- (i) Since $\det(AB) = \det(A) \cdot \det(B)$, it follows that if $\det(A) \neq 0$ and $\det(B) \neq 0$, then $\det(AB) \neq 0$, so $GL_n(\mathbf{F})$ is closed under matrix multiplication.
- (ii) Matrix multiplication is associative.
- (iii) $\det(A) \neq 0$ if and only if A has an inverse matrix, so each $A \in GL_n(\mathbf{F})$ has an inverse $A^{-1} \in GL_n(\mathbf{F})$ such that

$$AA^{-1} = A^{-1}A = I$$

where I is the $n \times n$ identity matrix.

Example 3.10 (Quaternion group). The *Quaternion group* Q_8 is defined by

$$Q_8 = \{1, -1, i, -i, j, -j, k, -k\}$$

with product \cdot computed as follows:

- $1 \cdot a = a \cdot 1 = a$ for all $a \in Q_8$

- $(-1) \cdot (-1) = 1$
- $(-1) \cdot a = a \cdot (-1) = -a$ for all $a \in Q_8$
- $i \cdot i = j \cdot j = k \cdot k = -1$
- $i \cdot j = k, j \cdot i = -k, j \cdot k = i, k \cdot j = -i, k \cdot i = j, i \cdot k = -j$

Note that Q_8 is a non-abelian group of order 8.

Subgroups

When given a set with certain properties, it is natural to consider its subsets that inherit the same properties.

Definition 3.11 (Subgroup). Let G be a group. Non-empty $H \subset G$ is a **subgroup** of G , denoted by $H \leq G$, if H is a group under the product in G .

That is, $H \leq G$ if and only if

- (i) $1 \in H$; (identity)
- (ii) $ab \in H$ for all $a, b \in H$; (closure)
- (iii) $a^{-1} \in H$ for all $a \in H$. (inverses)

Remark. There is no need to check that associativity holds in H , as it follows from associativity in G .

Every group G has two obvious subgroups: the group G itself, and the *trivial subgroup* $\{1\}$. A subgroup is a *proper subgroup* if it is not one of those two.

It would be useful to have some criterion for deciding whether a given subset of a group is a subgroup.

Lemma 3.12 (Subgroup criterion). Let G be a group. Then $H \leq G$ if and only if

- (i) $H \neq \emptyset$;
- (ii) $ab^{-1} \in H$ for all $a, b \in H$.

Proof.

\implies If $H \leq G$, then we are done, by definition of subgroup.

\impliedby Check group axioms:

- (i) Since $H \neq \emptyset$, there exists $a \in H$. Then $1 = aa^{-1} \in H$.
- (ii) Since $1 \in H$ and $a \in H$, then $a^{-1} = 1a^{-1} \in H$.
- (iii) For any $a, b \in H$, $a, b^{-1} \in H$, so by (ii), $a(b^{-1})^{-1} = ab \in H$.

□

Proposition 3.13. Let G be a group, $H, K \leq G$. Then $H \cap K \leq G$.

Proof. Apply the subgroup criterion:

- (i) Since $1 \in H$ and $1 \in K$, then $1 \in H \cap K$ so $H \cap K \neq \emptyset$.
- (ii) Let $a, b \in H \cap K$. Then $a, b \in H$ and $a, b \in K$. Since $H, K \leq G$, by the subgroup criterion, $ab^{-1} \in H$ and $ab^{-1} \in K$, so $ab^{-1} \in H \cap K$.

□

Corollary 3.14. Let G be a group, $\{H_i \mid i \in I\}$ is a collection of subgroups of G . Then

$$\bigcap_{i \in I} H_i \leq G.$$

Cyclic Groups

Definition 3.15. The *cyclic subgroup* H generated by $a \in G$, denoted by $H = \langle a \rangle$, is the set of all powers of a :

$$H = \{\dots, a^{-2}, a^{-1}, 1, a, a^2, \dots\}.$$

a is a *generator* of H .

You should verify that that $\langle a \rangle$ is indeed a subgroup of G . Furthermore, $\langle a \rangle$ is the smallest subgroup of G that contains a .

Remark. A cyclic subgroup may have more than one generator. For example, if $H = \langle a \rangle$, then also $H = \langle a^{-1} \rangle$ because $(a^{-1})^n = a^{-n} \in H$ for $n \in \mathbb{Z}$ so does $-n$, thus

$$\{a^n \mid n \in \mathbb{Z}\} = \{(a^{-1})^n \mid n \in \mathbb{Z}\}.$$

Lemma 3.16. Cyclic groups are abelian.

Proof. Let G be a cyclic group. For $a^i, a^j \in G$, by the laws of exponents,

$$a^i a^j = a^{i+j} = a^j a^i.$$

□

Proposition 3.17. A subgroup of a cyclic group is cyclic.

Proof. Let $a \in G$, $H \leq \langle a \rangle$. If $H = \{1\}$ then trivially H is cyclic.

Suppose that H contains some other element $b \neq 1$. Then $b = a^n$ for some integer n . Since H is a subgroup, $b^{-1} = a^{-n} \in H$. Since either n or $-n$ is positive, we can assume H contains positive powers of a and $n > 0$. Let m be the smallest positive integer such that $a^m \in H$ (such an m exist by the well-ordering principle).

Claim. $h = a^m$ is a generator for H .

We need to show that every $h' \in H$ can be written as a power of h . Since $h' \in H$ and $H \leq \langle a \rangle$, $h' = a^k$ for some integer k . By the division algorithm, there exist integers q, r such that $k = qm + r$ with $0 \leq r < m$. Hence

$$a^k = a^{qm+r} = (a^m)^q a^r = h^q a^r$$

so $a^r = a^k h^{-q}$. Since $a^k, h^{-q} \in H$, we must have $a^r \in H$. By the minimality of m , we must have $m = 0$ and so $k = qm$. Hence

$$h' = a^k = a^{qm} = h^q$$

and H is generated by h . □

Corollary 3.18. *The subgroups of \mathbb{Z} are exactly $n\mathbb{Z}$ for $n = 0, 1, 2, \dots$*

There are two possibilities: Either the powers x^n represent distinct elements, or they do not. We analyse the case that the powers of x are not distinct (finite cyclic groups).

Proposition 3.19. *Let $a \in G$, let $S = \{k \in \mathbb{Z} \mid a^k = 1\}$.*

(i) $S \leq \mathbb{Z}$.

(ii) $a^r = a^s$ (with $r \geq s$) if and only if $a^{r-s} = 1$.

(iii) *Suppose that S is not the trivial subgroup. Then $S = n\mathbb{Z}$ for some $n \in \mathbb{Z}^+$. The powers $1, a, a^2, \dots, a^{n-1}$ are the distinct elements of $\langle a \rangle$, and the order of $\langle a \rangle$ is n .*

Proof.

(i) $a^0 = 1$ so $0 \in S$.

If $k, l \in S$, $a^k = 1$ and $a^l = 1$ so $a^{k+l} = a^k a^l = 1$ then $k + l \in S$.

If $k \in S$, $a^k = 1$, then $a^{-k} = (a^k)^{-1} = 1$ so $-k \in S$.

(ii) This follows from the cancellation law.

(iii) Suppose that $S \neq \{0\}$. Then Corollary 3.18 shows that $S = n\mathbb{Z}$, where n is the smallest positive integer in S .

Let $a^k \in \langle a \rangle$. By the division algorithm, write $k = qn + r$ with $0 \leq r < n$. Then $a^{qn} = 1^q = 1$ so $a^k = a^{qn} a^r = a^r$. Hence a^k is equal to one of the powers $1, a, a^2, \dots, a^{n-1}$. It follows from (ii) that these powers are distinct, because a^n is the smallest positive power equal to 1. □

The group $\langle a \rangle = \{1, a, \dots, a^{n-1}\}$ described by (iii) above is called a *cyclic group of order n* .

More generally, we make the following definition.

Definition 3.20 (Subgroup generated by subset of group). Let G be a group, $S \subset G$. The *subgroup generated by S* , denoted by $\langle S \rangle$, is the smallest subgroup of G which contains S .

If $\langle S \rangle = G$, then the elements of S are said to be *generators* of G .

Notation. If $a \in G$, then we write $\langle a \rangle$ (rather than the more accurate but cumbersome $\langle \{a\} \rangle$).

Order

Definition 3.21 (Order). Let G be a group, $a \in G$. If there is a positive integer k such that $a^k = 1$, then the **order** of a is defined as

$$o(a) := \min\{m > 0 \mid a^m = 1\}.$$

Otherwise we say that the order of a is infinite.

Proposition 3.22. *If G is finite, then $o(a)$ is finite for each $a \in G$.*

Proof. Consider the list

$$a, a^2, a^3, \dots \in G.$$

Since G is finite, this list must have repeats. Hence $a^i = a^j$ for some integers $i > j$, so $a^{i-j} = 1$. This shows that $\{m > 0 \mid a^m = e\}$ is non-empty and thus has a minimal element. \square

Proposition 3.23. *If $a \in G$ and $o(a)$ is finite, then $a^n = 1$ if and only if $o(a) \mid n$.*

Proof.

\Leftarrow Suppose $o(a) \mid n$. Then $n = ko(a)$ for some $k \in \mathbb{Z}$, so

$$a^n = \left(a^{o(a)}\right)^k = 1^k = 1.$$

\Rightarrow Suppose $a^n = 1$. By the division algorithm, there exists integers q, r such that $n = qo(a) + r$, where $0 \leq r < o(a)$. Then

$$a^r = a^{n-qo(a)} = a^n \left(a^{o(a)}\right)^{-q} = 1.$$

By the minimality of $o(a)$, we must have $r = 0$, and so $n = qo(a)$ implies $o(a) \mid n$. \square

Corollary 3.24. *Let G be a cyclic group, $a \in G$. Then $a^k = a^m$ if and only if $m \equiv k \pmod{o(a)}$.*

§3.2 Cosets

Definition 3.25 (Coset). Let $H \leq G$. For $a \in G$, a *left coset* and *right coset* of H in G are

$$aH := \{ah \mid h \in H\}$$

$$Ha := \{ha \mid h \in H\}$$

Any element of a coset is called a *representative* for the coset.

The set of left cosets is given by

$$(G/H)_l := \{aH \mid a \in G\}.$$

Similarly, the set of right cosets is given by

$$(G/H)_r := \{Ha \mid a \in G\}.$$

Lemma 3.26. Let $H \leq G$. Then $aH = H$ if and only if $a \in H$. (Similarly, $Ha = H$ if and only if $a \in H$.)

Proof.

\implies Suppose $aH = H$. Then $ah \in H$ for some $h \in H$. Let $k = ah$, then $a = kh^{-1} \in H$.

\impliedby Let $a \in H$. Then $aH \subset H$.

Since $a^{-1} \in H$, $a^{-1}H \subset H$. Then $H = eH = (aa^{-1})H = a(a^{-1})H \subset aH$. Hence $aH = H$. \square

The next result shows when two cosets are equal.

Lemma 3.27. Let $H \leq G$, $a, b \in G$. Then $aH = bH$ if and only if $a^{-1}b \in H$.

Proof.

$$\begin{aligned} aH = bH &\iff a^{-1}(aH) = a^{-1}bH \\ &\iff (a^{-1}a)H = (a^{-1}b)H \\ &\iff H = (a^{-1}b)H \end{aligned}$$

Note that from the previous result, $H = (a^{-1}b)H$ if and only if $a^{-1}b \in H$. \square

Proposition 3.28. Let $H \leq G$. Then $(G/H)_l$ forms a partition of G . (Similar remarks hold for right cosets.)

We need to prove the following.

- (i) For all $a \in G$, $aH \neq \emptyset$.
- (ii) $\bigcup_{a \in G} aH = G$.
- (iii) For every $a, b \in G$, $aH \cap bH = \emptyset$ or $aH = bH$.

Proof.

- (i) Since $H \leq G$, $e \in H$. Thus for all $a \in G$, $a = ae \in aH$ so $aH \neq \emptyset$.
- (ii) For all $a \in G$, $aH \subset G$, then $\bigcup_{a \in G} aH \subset G$. Note that $a \in G$ implies $a = ae \in aH$, and so $G = \bigcup_{a \in G} aH \subset \bigcup_{a \in G} aH$. By double inclusion we are done.
- (iii) If $aH \cap bH = \emptyset$, then we are done. If $aH \cap bH \neq \emptyset$ we need to show $aH = bH$. Let $x \in G$ such that $x \in aH \cap bH$. Then $x = ah_1 = bh_2$ for $h_1, h_2 \in H$ so $h_1 = a^{-1}bh_2$. Notice that $a^{-1}b = h_1h_2^{-1} \in H$ and thus $aH = bH$.

□

Lagrange's Theorem

Definition 3.29 (Index). Let $H \leq G$. The *index* of H in G is the number of left cosets of H in G , denoted by $|G : H|$.

The following result shows that H partitions G into equal-sized parts.

Lemma 3.30. *The cosets of H in G are the same size as H ; that is, for all $a \in G$, $|aH| = |H|$.*

Proof. Let $f: H \rightarrow aH$ which sends $h \mapsto ah$. For $h_1, h_2 \in H$,

$$\begin{aligned} f(h_1) = f(h_2) &\implies ah_1 = ah_2 \\ &\implies a^{-1}ah_1 = a^{-1}ah_2 \\ &\implies h_1 = h_2 \end{aligned}$$

thus f is an injective mapping. Note that f is surjective by the definition of aH . Since f is bijective, $|H| = |aH|$. □

Theorem 3.31 (Lagrange's theorem). *Let G be a finite group, $H \leq G$. Then $|G| = |H| |G : H|$.*

Proof. Let $|H| = n$, and let $|G : H| = k$. Since G is partitioned into k disjoint subsets, each of which has cardinality n , we have $|G| = kn$, or

$$|G| = |H| |G : H| \tag{3.1}$$

as desired. □

Eq. (3.1) is known as the *counting formula*.

Corollary 3.32. *The order of an element of a finite group divides the order of the group.*

Proof. Let $a \in G$. Then by Proposition 3.19, $o(a) = |\langle a \rangle|$.

Since $\langle a \rangle$ is a subgroup of G , by Lagrange's Theorem, $|\langle a \rangle|$ divides $|G|$; that is, $o(a)$ divides $|G|$. □

Corollary 3.33. *A group of prime order is cyclic.*

Proof. Let $|G| = p$ be prime. Let $a \in G, a \neq 1$. We will show that $G = \langle a \rangle$.

Since $o(a) \mid |G| = p$ and $o(a) > 1$, we must have $o(a) = p$. Notice that this is also the order of $\langle a \rangle$. Since G has order p , thus $\langle a \rangle = G$. \square

This corollary classifies groups of prime order p . They form one isomorphism class: the class of the cyclic groups of order p .

The next result is of great interest in number theory. The *Euler ϕ -function* $\phi(n)$ is defined for all positive integers as follows:

$$\phi(n) = \begin{cases} 1 & (n = 1) \\ \text{number of positive integers less than } n, \text{ relatively prime to } n & (n > 1) \end{cases}$$

Theorem 3.34 (Euler). *If n is a positive integer, and a is coprime to n , then*

$$a^{\phi(n)} \equiv 1 \pmod{n}.$$

Theorem 3.35 (Fermat). *If p is prime, and a is any integer, then*

$$a^p \equiv a \pmod{p}.$$

Proof. If n is a prime number p , then $\phi(p) = p - 1$. We consider two cases.

- If a is coprime to p , then by Euler's totient theorem, $a^{p-1} \equiv 1 \pmod{p}$, and the desired result follows.
- If a is not coprime to p , since p is prime, we must have that $p \mid a$, so that $a \equiv 0 \pmod{p}$. Hence $0 \equiv a^p \equiv a \pmod{p}$ here also.

\square

Counting Principle

We generalise the notion of cosets, as defined earlier.

Definition 3.36. Let $H, K \leq G$, define

$$HK = \{hk \mid h \in H, k \in K\}.$$

Lemma 3.37. *$HK \leq G$ if and only if $HK = KH$.*

Proof.

\Leftarrow Suppose $HK = KH$; that is, if $h \in H$ and $k \in K$, then $hk = k_1h_1$ for some $k_1 \in K, h_1 \in H$.

We now show that HK is a subgroup of G :

(i) $1 \in H$ and $1 \in K$, so $1 \in HK$.

(ii) Let $x = hk \in HK$, $y = h'k' \in HK$. then

$$xy = hkh'k'.$$

Note that $kh' \in KH = HK$, so $kh' = h_2k_2$ for some $h_2 \in H, k_2 \in K$. Then

$$xy = h(h_2k_2)k' = (hh_2)(k_2k') \in HK.$$

Thus HK is closed.

(iii) Let $x \in HK$, then $x = hk$ for some $h \in H, k \in K$. Thus

$$x^{-1} = (hk)^{-1} = k^{-1}h^{-1} \in KH = HK,$$

so $x^{-1} \in HK$.

\implies Suppose $HK \leq G$.

• Let $x \in KH$, so $x = kh$ for some $k \in K, h \in H$. Then

$$x = kh = (h^{-1}k^{-1})^{-1} \in HK.$$

Thus $KH \subset HK$.

• Let $x \in HK$. Since $HK \leq G$, HK is closed under inverses, so $x^{-1} = hk \in HK$. Then

$$x = (x^{-1})^{-1} = (hk)^{-1} = k^{-1}h^{-1} \in KH.$$

Thus $HK \subset KH$.

Hence $HK = KH$. □

An interesting special case is the situation when G is an abelian group, for in that case trivially $HK = KH$. Thus as a consequence we have the following result.

Corollary 3.38. *Let $H, K \leq G$, where G is abelian. Then $HK \leq G$.*

Proposition 3.39. *If $H, K \leq G$ are finite groups, then*

$$|HK| = \frac{|H||K|}{|H \cap K|}.$$

Proof. Notice that HK is a union of left cosets of K , namely

$$HK = \bigcup_{h \in H} hK.$$

□

Normal Subgroups, Quotient Groups

Definition 3.40 (Normal subgroup). Let G be a group. $H \leq G$ is a **normal subgroup** of G , denoted by $H \triangleleft G$, if

$$aH = Ha \quad (\forall a \in G)$$

If G has no non-trivial normal subgroup, then G is a *simple group*.

Remark. This does *not* mean that $ah = ha$ for all $a \in G, h \in H$ or that G is abelian. Although we can easily see that all subgroups of abelian groups are normal. In general, a left coset does not equal the right coset.

Lemma 3.41. *The following are equivalent.*

- (i) $H \triangleleft G$.
- (ii) $ghg^{-1} \in H$ for all $g \in G, h \in H$.
- (iii) $gHg^{-1} = H$ for all $g \in G$.

Proof.

(i) \iff (ii) In the forward direction, $aH = Ha$ for all $a \in G$. Let $g \in G, x \in H$. Then $gH = Hg$ so $gx = h'g$ for some $h' \in H$. Then $gHg^{-1} = h'gg^{-1} = h' \in H$.

In the reverse direction, $ghg^{-1} \in H$ for all $g \in G, h \in H$. Fix g . Then $ghg^{-1} \in H$ implies $gh \in Hg$ for all $h \in H$. So $gH \subset Hg$. Similarly $gH \supset Hg$, so $gH = Hg$.

(i) \iff (iii) $H \triangleleft G$ if and only if for all $g \in G$,

$$\begin{aligned} gH = Hg &\iff (gH)g^{-1} = (Hg)g^{-1} \\ &\iff gHg^{-1} = H \end{aligned}$$

□

Remark. We frequently use (ii) to check if a subgroup is a normal subgroup.

Definition 3.42 (Quotient group). Let G be a group, $H \triangleleft G$. Then the **quotient group** of G by H is

$$G/H := \{aH \mid a \in G\}.$$

Lemma 3.43. G/H is a group under the following operation: for all $aH, bH \in G/H$,

$$(aH)(bH) = a(Hb)H = a(bH)H = abH$$

Proof. Check group axioms.

- (i) For $a, b, c \in G$,

$$(aH)(bHcH) = (aH)(bcH) = a(bc)H = (ab)cH = (aHbH)cH$$

so the operation is associative.

(ii) The identity of G/H is the coset $eH = H$.

(iii) For $aH \in G/H$, the inverse of aH is $a^{-1}H$ as is immediate from the definition of the product:

$$(aH)(a^{-1}H) = aa^{-1}H = H \implies (aH)^{-1} = a^{-1}H.$$

□

Lemma 3.44. *Let G be a finite group, $H \triangleleft G$. Then*

$$|G/H| = |G : H| = \frac{|G|}{|H|}.$$

Proof. Since G/H has as its elements the left cosets of H in G , and since there are precisely $|G : H|$ such cosets, by Lagrange's theorem, we obtain the desired result. □

Definition 3.45 (Quotient map). Let $H \triangleleft G$. The *quotient map* is the map $\pi : G \rightarrow G/H$ which sends $a \mapsto aH$.

§3.3 Homomorphisms and Isomorphisms

In this section, we make precise the notion of when two groups “look the same”; that is, they have the same group-theoretic structure. This is the notion of an *isomorphism* between two groups.

When we talk about functions between groups it makes sense to limit our scope to functions that preserve the group operation (morphisms in the category of groups). More precisely:

Definition 3.46 (Homomorphism). Let $(G, *)$ and (H, \diamond) be groups. A map $\phi: G \rightarrow H$ is called a *homomorphism* if, for all $x, y \in G$,

$$\phi(x * y) = \phi(x) \diamond \phi(y).$$

When the group operations for G and H are not explicitly written, we have

$$\phi(xy) = \phi(x)\phi(y).$$

Definition 3.47 (Isomorphism). An *isomorphism* $\phi: G \rightarrow H$ is a bijective homomorphism. If $\phi: G \rightarrow H$ is an isomorphism, then G and H are *isomorphic*, denoted by $G \cong H$.

An *automorphism* of a group G is an isomorphism from G to G ; the automorphisms of G form a group $\text{Aut}(G)$ under composition. An *endomorphism* of G is a homomorphism from G to G .

Example 3.48. $(\mathbb{R}, +) \cong (\mathbb{R}^+, \times)$, as the exponential map $\exp: \mathbb{R} \rightarrow \mathbb{R}^+$ defined by $\exp(x) = e^x$ is an isomorphism from $(\mathbb{R}, +)$ to (\mathbb{R}^+, \times) .

- (i) \exp is a bijection since it has an inverse function (namely \ln).
- (ii) \exp preserves the group operations since $e^{x+y} = e^x e^y$.

Proposition 3.49. Let $\phi: G \rightarrow H$ be a homomorphism. Let $g \in G$, $n \in \mathbb{Z}$. Then

- (i) $\phi(1_G) = 1_H$;
- (ii) $\phi(g^{-1}) = (\phi(g))^{-1}$;
- (iii) $\phi(g^n) = (\phi(g))^n$.

Proof.

- (i) $\phi(1_G) = \phi(1_G 1_G) = \phi(1_G)\phi(1_G)$, then apply $\phi(1_G)^{-1}$ to both sides to get $\phi(1_G) = 1_H$.
- (ii) $\phi(g)\phi(g^{-1}) = \phi(gg^{-1}) = \phi(1_G) = 1_H$.
- (iii) Note more generally that we can show $\phi(g^n) = (\phi(g))^n$ for $n > 0$ by induction. For $n = -k < 0$ we have

$$\phi(g^n) = \phi((g^{-1})^k) = (\phi(g^{-1}))^k = (\phi(g)^{-1})^k = \phi(g)^n.$$

□

Proposition 3.50. *Quotient maps are homomorphisms.*

Proof. Let $\pi: G \rightarrow G/H$ which sends $g \mapsto gH$ be a quotient map. Then for all $x, y \in G$,

$$\pi(xy) = (xy)H = (xH)(yH) = \pi(x)\pi(y).$$

□

Kernel and Image

Definition 3.51 (Kernel and image). Let $\phi: G \rightarrow H$ be a homomorphism. Then the *kernel* of ϕ is

$$\ker \phi := \{g \in G \mid \phi(g) = 1_H\} \subset G.$$

The *image* of G under ϕ is

$$\text{im } \phi := \phi(G) = \{\phi(g) \mid g \in G\} \subset H.$$

Remark. $\text{im } \phi$ is the usual set theoretic image of ϕ .

Proposition 3.52. *Let $\phi: G \rightarrow H$ be a homomorphism. Then*

(i) $\ker \phi \triangleleft G$;

(ii) $\text{im } \phi \leq H$.

Proof.

(i) Apply the subgroup criterion. Since $1_G \in \ker \phi$, $\ker \phi \neq \emptyset$. Let $x, y \in \ker \phi$; that is, $\phi(x) = \phi(y) = 1_H$. Then

$$\phi(xy^{-1}) = \phi(x)\phi(y)^{-1} = 1_H$$

so $xy^{-1} \in \ker \phi$. By the subgroup criterion, $\ker \phi \leq G$.

Let $x \in \ker \phi$, $g \in G$. Then

$$\phi(gxg^{-1}) = \phi(g)\phi(x)\phi(g^{-1}) = 1,$$

so $gxg^{-1} \in \ker \phi$. Hence $\ker \phi \triangleleft G$.

(ii) Since $\phi(1_G) = 1_H$, $1_H \in \text{im } \phi$ so $\text{im } \phi \neq \emptyset$. Let $x, y \in \text{im } \phi$. Then there exists $a, b \in G$ such that $\phi(a) = x$, $\phi(b) = y$. Then

$$xy^{-1} = \phi(a)\phi(b)^{-1} = \phi(ab^{-1})$$

so $xy^{-1} \in \text{im } \phi$. By the subgroup criterion, $\text{im } \phi \leq G$.

□

The following result is a useful characterisation for injective homomorphisms.

Lemma 3.53. *Let $\phi: G \rightarrow H$ be a homomorphism. Then ϕ is injective if and only if $\ker \phi = \{1_G\}$.*

Proof.

\Rightarrow Suppose ϕ is injective. Since $\phi(1_G) = 1_H$, $1_G \in \ker \phi$ so $\{1_G\} \subset \ker \phi$.

Conversely, let $x \in \ker \phi$, so $\phi(x) = 1_H$. Then $\phi(x) = 1_H = \phi(1_G)$, so by injectivity $x = 1_G$. Hence $\ker \phi \subset \{1_G\}$, so $\ker \phi = \{1_G\}$.

\Leftarrow Suppose $\ker \phi = \{1_G\}$. Suppose $\phi(a) = \phi(b)$, then $\phi(ab^{-1}) = \phi(a)\phi(b^{-1}) = \phi(a)\phi(a)^{-1} = 1_H$. Hence $ab^{-1} \in \ker \phi = \{1_G\}$, so $ab^{-1} = 1_G$ and thus $a = b$. Therefore ϕ is injective. \square

Lemma 3.54. *Let $\phi: G \rightarrow H$ be an isomorphism. Then the inverse map $\phi^{-1}: H \rightarrow G$ is also an isomorphism.*

Proof. The inverse of a bijective map is bijective. Hence it suffices to show that $\phi^{-1}(x)\phi^{-1}(y) = \phi^{-1}(xy)$ for all $x, y \in H$.

Let $a = \phi^{-1}(x)$, $b = \phi^{-1}(y)$, $c = \phi^{-1}(xy)$; we will show that $ab = c$. Since ϕ is bijective, it suffices to show that $\phi(ab) = \phi(c)$.

Since ϕ is a homomorphism,

$$\phi(ab) = \phi(a)\phi(b) = xy = \phi(c).$$

\square

Isomorphism Theorems

Theorem 3.55 (First isomorphism theorem). *Let $\phi: G \rightarrow H$ be a homomorphism. Then*

$$G / \ker \phi \cong \text{im } \phi(G).$$

Proof. Let $K = \ker \phi$. Let

$$\begin{aligned} \theta: G/K &\rightarrow \text{im } \phi \\ \forall x \in G, \quad xK &\mapsto \phi(x) \end{aligned}$$

Claim. θ is an isomorphism.

We first need to check if θ is well-defined: let $x, y \in G$. Suppose $xK = yK$. Then

$$\begin{aligned} xK = yK & \\ \iff x^{-1}y \in K & \\ \iff \phi(x^{-1}y) = 1_H & \\ \iff \phi(x)^{-1}\phi(y) = 1_H & \\ \iff \phi(x) = \phi(y) & \end{aligned}$$

Next we show that θ is a homomorphism: for all $x, y \in G$,

$$\theta(xKyK) = \theta(xyK) = \phi(xy) = \phi(x)\phi(y) = \theta(xK)\theta(yK).$$

Finally we show that θ is bijective:

- θ is injective since

$$\theta(xK) = \theta(yK) \implies \phi(x) = \phi(y) \implies xK = yK.$$

- θ is surjective, since

$$\begin{aligned} \text{im } \theta &= \{\theta(xK) \mid x \in G\} \\ &= \{\phi(x) \mid x \in G\} \\ &= \text{im } \phi. \end{aligned}$$

□

Theorem 3.56 (Second isomorphism theorem). *Let $A \leq G$, $B \triangleleft G$. Then*

(i) $AB \leq G$;

(ii) $B \triangleleft AB$;

(iii) $A \cap B \triangleleft A$;

(iv) $AB/B \cong A/(A \cap B)$.

Theorem 3.57 (Third isomorphism theorem). *Let $H, K \triangleleft G$, $H \leq K$. Then $K/H \triangleleft G/H$, and*

$$(G/H)/(K/H) \cong G/K.$$

If we denote the quotient by H with a bar, this can be written

$$\overline{G}/\overline{K} \cong G/K.$$

Theorem 3.58 (Fourth isomorphism theorem).

Theorem 3.59 (Cayley's theorem).

§3.4 Group Actions

We move now, from thinking of groups in their own right, to thinking of how groups can move sets around—for example, how S_n permutes $\{1, 2, \dots, n\}$ and matrix groups move vectors.

Definition 3.60 (Group action). A **group action** of a group G on a set A is a map from $G \times A \rightarrow A$ (written as $g \cdot a$, for all $g \in G, a \in A$) satisfying the following properties:

- (i) $g_1 \cdot (g_2 \cdot a) = (g_1 g_2) \cdot a$, for all $g_1, g_2 \in G, a \in A$;
- (ii) $1_G \cdot a = a$ for all $a \in A$.

We say that G is a group acting on a set A .

Intuitively, a group action of G on a set A means that every element g in G acts as a permutation on A in a manner consistent with the group operations in G . There is also a notion of *left action* and *right action*.

For the following definitions, let G be a group, and $A \subset G$ be non-empty.

Definition 3.61 (Centraliser). The **centraliser** of A in G is defined by

$$C_G(A) := \{g \in G \mid \forall a \in A, gag^{-1} = a\}.$$

Since $gag^{-1} = a$ if and only if $ga = ag$, $C_G(A)$ is the set of elements of G which commute with every element of A .

We check that $C_G(A) \leq G$:

- (i) $e \in C_G(A)$, so $C_G(A) \neq \emptyset$.
- (ii) Let $x, y \in C_G(A)$; that is, for all $a \in A$, $xax^{-1} = a$ and $yay^{-1} = a$. Then

$$\begin{aligned} (xy)a(xy)^{-1} &= (xy)a(y^{-1}x^{-1}) \\ &= x(yay^{-1})x^{-1} \\ &= xax^{-1} = a \end{aligned}$$

so $xy \in C_G(A)$. Hence $C_G(A)$ is closed under products.

- (iii) Let $x \in C_G(A)$; that is, for all $a \in A$, $xax^{-1} = a$. Applying x^{-1} to both sides gives $ax^{-1} = x^{-1}a$. Applying x to both sides gives $a = x^{-1}ax$, so $x^{-1} \in C_G(A)$. Hence $C_G(A)$ is closed under taking inverses.

Notation. In the special case when $A = \{a\}$ we simply write $C_G(a)$ instead of $C_G(\{a\})$. In this case $a^n \in C_G(a)$ for all $n \in \mathbb{Z}$.

Definition 3.62 (Centre). The **centre** of G is the set of elements which commute with all the elements of G :

$$Z(G) := \{g \in G \mid \forall x \in G, gx = xg\}.$$

Note that $Z(G) = C_G(G)$, so the argument above proves $Z(G) \leq G$ as a special case.

Definition 3.63 (Normaliser). Define $gAg^{-1} = \{gag^{-1} \mid a \in A\}$. The *normaliser* of A in G is

$$N_G(A) := \{g \in G \mid gAg^{-1} = A\}.$$

Notice that if $g \in C_G(A)$, then $gag^{-1} = a \in A$ for all $a \in A$ so $C_G(A) \leq N_G(A)$. The proof that $N_G(A) \leq G$ is similar to the one that $C_G(A) \leq G$.

Definition 3.64 (Stabiliser). If G is a group acting on a set S , $s \in S$, then the *stabiliser* of s in G is

$$G_s := \{g \in G \mid g \cdot s = s\}.$$

Notation. Denote the set of all fixed points to be $S^G = \{s \in S \mid \forall g \in G, gs = g\}$.

We check that $G_s \leq G$:

(i) By definition of group action, $1_G \cdot a = a$, so $1_G \in G_s$.

(ii) Let $x, y \in G_s$, then

$$\begin{aligned} (xy) \cdot s &= x \cdot (y \cdot s) \\ &= x \cdot s = s \end{aligned}$$

so $xy \in G_s$. Hence G_s is closed under products.

(iii) Let $x \in G_s$; that is, $x \cdot s = s$. Then

$$\begin{aligned} x^{-1} \cdot s &= x^{-1} \cdot (x \cdot s) \\ &= (x^{-1}x) \cdot s \\ &= e \cdot s = s \end{aligned}$$

so $x^{-1} \in G_s$. Hence G_s is closed under taking inverses.

Definition 3.65. The *kernel* of the action of G on S is

$$\{g \in G \mid \forall s \in S, g \cdot s = s\}.$$

Definition 3.66 (Orbit). Let G be a group that acts on a set S . Define the *orbit* of a group element $s \in S$ as

$$G(s) := \{g \cdot s \in S \mid g \in G\}.$$

Conjugation

Sylow's Theorem

Definition 3.67 (Sylow p -subgroup). Let G be a group, and let p be a prime.

- (i) A group of order p^α ($\alpha \geq 1$) is called a p -group. Subgroups of G which are p -groups are called p -subgroups.
- (ii) If $|G| = p^\alpha m$ ($p \nmid m$), then a subgroup of order p^α is called a **Sylow p -subgroup** of G .

Notation. The set of Sylow p -subgroups of G is denoted by $Syl_p(G)$, and the number of Sylow p -subgroups of G is denoted by $n_p(G)$ (or just n_p when G is clear from the context).

Theorem 3.68 (Sylow's theorem). Let $|G| = p^\alpha m$, where p is a prime and $p \nmid m$.

- (i) Sylow p -subgroups of G exist, i.e. $Syl_p(G) \neq \emptyset$.
- (ii) If P is a Sylow p -subgroup of G , and Q is any p -subgroup of G , then there exists $g \in G$ such that $Q \leq gPg^{-1}$, i.e. Q is contained in some conjugate of P . In particular, any two Sylow p -subgroups of G are conjugate in G .
- (iii) $n_p \equiv 1 \pmod{p}$. Furthermore, n_p is the index in G of the normaliser $N_G(P)$ for any Sylow p -subgroup P , hence $n_p \mid m$.

§3.5 Group Product, Finite Abelian Groups

Definition 3.69 (Direct product). The *direct product* $G_1 \times \cdots \times G_n$ of the groups $(G_1, *_1), \dots, (G_n, *_n)$ is the Cartesian product

$$G_1 \times \cdots \times G_n := \{(g_1, \dots, g_n) \mid g_i \in G_i\}$$

with operation defined componentwise:

$$(g_1, \dots, g_n) * (h_1, \dots, h_n) = (g_1 *_1 h_1, \dots, g_n *_n h_n).$$

Proposition 3.70. *If G_1, \dots, G_n are groups, then*

$$|G_1 \times \cdots \times G_n| = |G_1| |G_2| \cdots |G_n|.$$

Proof. Let $G = G_1 \times \cdots \times G_n$. The proof that the group axioms hold for G is straightforward since each axiom is a consequence of the fact that the same axiom holds for each G_i , and the operation on G defined componentwise.

The number of n -tuples in G follows from simple combinatorics. □

Exercises

Exercise 3.1. Show that any two cyclic groups of the same order are isomorphic.

Solution. Suppose $\langle x \rangle$ and $\langle y \rangle$ are both cyclic groups of order n . We first prove the case where $n < \infty$. We claim that the map $\phi: \langle x \rangle \rightarrow \langle y \rangle$ which sends $x^k \mapsto y^k$ is an isomorphism.

Lemma. Let G be a group, $g \in G$, let $m, n \in \mathbb{Z}$. Denote $d = \gcd(m, n)$. If $g^n = 1$ and $g^m = 1$, then $g^d = 1$.

Proof. By Bezout's lemma, since $d = \gcd(m, n)$, then there exists $q, r \in \mathbb{Z}$ such that $qm + rn = d$. Thus

$$g^d = g^{qm+rn} = (g^m)^q (g^n)^r = 1.$$

□

We first show that ϕ is well-defined; that is, $x^r = x^s \implies \phi(x^r) = \phi(x^s)$. Note that $x^{r-s} = e$, so by the above lemma, $n \mid r - s$. Write $r = tn + s$ for some $t \in \mathbb{Z}$, so

$$\phi(x^r) = \phi(x^{tn+s}) = y^{tn+s} = (y^n)^t y^s = y^s = \phi(x^s).$$

We then show that ϕ is a homomorphism:

$$\phi(x^a x^b) = \phi(x^{a+b}) = y^{a+b} = y^a y^b = \phi(x^a) \phi(x^b).$$

Finally we show that ϕ is bijective. Since the element y^k of $\langle y \rangle$ is in the image of x^k under ϕ , ϕ is surjective. Since both groups have the same finite order, any surjection from one to the other is a bijection. Therefore ϕ is an isomorphism.

We now prove the case where $n = \infty$. If $\langle x \rangle$ is an infinite cyclic group, let $\phi: \mathbb{Z} \rightarrow \langle x \rangle$ be defined by $\phi(k) = x^k$. (This map is well-defined since there is no ambiguity in the representation of elements in the domain.)

Since $x^a \neq x^b$ for all distinct $a, b \in \mathbb{Z}$, ϕ is injective. By definition of a cyclic group, ϕ is surjective. As above, the laws of exponents ensure ϕ is a homomorphism. Hence ϕ is an isomorphism. □

III

Linear Algebra

4 Finite Dimensional Vector Spaces

§4.1 Definition of Vector Space

Notation. A field is denoted by \mathbf{F} , which can mean either \mathbb{R} or \mathbb{C} . \mathbf{F}^n is the set of n -tuples whose elements belong to \mathbf{F} :

$$\mathbf{F}^n := \{(x_1, \dots, x_n) \mid x_i \in \mathbf{F}\}$$

For $(x_1, \dots, x_n) \in \mathbf{F}^n$ and $i = 1, \dots, n$, we say that x_i is the i -th coordinate of (x_1, \dots, x_n) .

Definition 4.1 (Vector space). V is a **vector space** over \mathbf{F} if the following properties hold:

(i) Addition is commutative: $u + v = v + u$ for all $u, v \in V$

(ii) Addition is associative: $(u + v) + w = u + (v + w)$ for all $u, v, w \in V$

Multiplication is associative: $(ab)v = a(bv)$ for all $v \in V, a, b \in \mathbf{F}$

(iii) Additive identity: there exists $\mathbf{0} \in V$ such that $v + \mathbf{0} = v$ for all $v \in V$

(iv) Additive inverse: for every $v \in V$, there exists $w \in V$ such that $v + w = \mathbf{0}$

(v) Multiplicative identity: $1v = v$ for all $v \in V$

(vi) Distributive properties: $a(u + v) = au + av$ and $(a + b)v = av + bv$ for all $a, b \in \mathbf{F}$ and $u, v \in V$

Notation. For the rest of this text, V denotes a vector space over \mathbf{F} .

Example 4.2. \mathbb{R}^n is a vector space over \mathbb{R} , \mathbb{C}^n is a vector space over \mathbb{C} .

Elements of a vector space are called *vectors* or *points*.

The scalar multiplication in a vector space depends on \mathbf{F} . Thus when we need to be precise, we will say that V is a vector space over \mathbf{F} instead of saying simply that V is a vector space. For example, \mathbb{R}^n is a vector space over \mathbb{R} , and \mathbb{C}^n is a vector space over \mathbb{C} . A vector space over \mathbb{R} is called a *real vector space*; a vector space over \mathbb{C} is called a *complex vector space*.

Lemma 4.3 (Uniqueness of additive identity). *A vector space has a unique additive identity.*

Proof. Suppose otherwise, then $\mathbf{0}$ and $\mathbf{0}'$ are additive identities of V . Then

$$\mathbf{0}' = \mathbf{0}' + \mathbf{0} = \mathbf{0} + \mathbf{0}' = \mathbf{0}$$

where the first equality holds because $\mathbf{0}$ is an additive identity, the second equality comes from commutativity, and the third equality holds because $\mathbf{0}'$ is an additive identity. Thus $\mathbf{0}' = \mathbf{0}$. \square

Lemma 4.4 (Uniqueness of additive inverse). *Every element in a vector space has a unique additive inverse.*

Proof. Suppose otherwise, then for $v \in V$, w and w' are additive inverses of v . Then

$$w = w + \mathbf{0} = w + (v + w') = (w + v) + w' = \mathbf{0} + w' = w'.$$

Thus $w = w'$. □

Because additive inverses are unique, the following notation now makes sense.

Notation. Let $v, w \in V$. Then $-v$ denotes the additive inverse of v ; $w - v$ is defined to be $w + (-v)$.

We now prove some seemingly trivial facts.

Lemma 4.5.

(i) For every $v \in V$, $0v = \mathbf{0}$.

(ii) For every $a \in \mathbf{F}$, $a\mathbf{0} = \mathbf{0}$.

(iii) For every $v \in V$, $(-1)v = -v$.

Proof.

(i) Let $v \in V$,

$$0v = (0 + 0)v = 0v + 0v.$$

Adding the additive inverse of $0v$ to both sides of the equation gives $\mathbf{0} = 0v$.

(ii) Let $a \in \mathbf{F}$,

$$a\mathbf{0} = a(\mathbf{0} + \mathbf{0}) = a\mathbf{0} + a\mathbf{0}.$$

Adding the additive inverse of $a\mathbf{0}$ to both sides of the equation gives $\mathbf{0} = a\mathbf{0}$.

(iii) Let $v \in V$,

$$v + (-1)v = 1v + (-1)v = (1 + (-1))v = 0v = \mathbf{0}.$$

Since $v + (-1)v = \mathbf{0}$, $(-1)v$ is the additive inverse of v . □

Example 4.6. \mathbf{F}^∞ is defined to be the set of all sequences of elements of \mathbf{F} :

$$\mathbf{F}^\infty := \{(x_1, x_2, \dots) \mid x_i \in \mathbf{F}\}$$

Define addition and scalar multiplication on \mathbf{F}^∞ as

$$\begin{aligned} (x_1, x_2, \dots) + (y_1, y_2, \dots) &= (x_1 + y_1, x_2 + y_2, \dots) \\ \lambda(x_1, x_2, \dots) &= (\lambda x_1, \lambda x_2, \dots) \end{aligned}$$

Then \mathbf{F}^∞ is a vector space over \mathbf{F} , where the additive identity is $\mathbf{0} = (0, 0, \dots)$.

Our next example of a vector space involves a set of functions.

Example 4.7. If S is a set, $\mathbf{F}^S := \{f \mid f : S \rightarrow \mathbf{F}\}$. Define addition and scalar multiplication on \mathbf{F}^S as

$$\begin{aligned}(f + g)(x) &= f(x) + g(x) \quad (x \in S) \\ (\lambda f)(x) &= \lambda f(x) \quad (x \in S)\end{aligned}$$

for all $f, g \in \mathbf{F}^S$, $\lambda \in \mathbf{F}$. Then \mathbf{F}^S is a vector space over \mathbf{F} (if S is a non-empty set), where the additive identity of \mathbf{F}^S is the function $0 : S \rightarrow \mathbf{F}$ defined as

$$0(x) = 0 \quad (\forall x \in S)$$

and for $f \in \mathbf{F}^S$, additive inverse of f is the function $-f : S \rightarrow \mathbf{F}$ defined as

$$(-f)(x) = -f(x) \quad (\forall x \in S)$$

Remark. \mathbf{F}^n and \mathbf{F}^∞ are special cases of the vector space \mathbf{F}^S ; think of \mathbf{F}^n as $\mathbf{F}^{\{1,2,\dots,n\}}$, and \mathbf{F}^∞ as $\mathbf{F}^{\{1,2,\dots\}}$.

Example 4.8 (Complexification). Suppose V is a real vector space. The *complexification* of V , denoted by $V_{\mathbb{C}}$, equals $V \times V$. An element of $V_{\mathbb{C}}$ is an ordered pair (u, v) , where $u, v \in V$, which we write as $u + iv$.

- Addition on $V_{\mathbb{C}}$ is defined as

$$(u_1 + iv_1) + (u_2 + iv_2) = (u_1 + u_2) + i(v_1 + v_2)$$

for all $u_1, v_1, u_2, v_2 \in V$.

- Complex scalar multiplication on $V_{\mathbb{C}}$ is defined as

$$(a + bi)(u + iv) = (au - bv) + i(av + bu)$$

for all $a, b \in \mathbb{R}$ and all $u, v \in V$.

Then $V_{\mathbb{C}}$ is a (complex) vector space.

§4.2 Subspaces

Whenever we have a mathematical object with some structure, we want to consider subsets that also have the same structure.

Definition 4.9 (Subspace). $U \subset V$ is a **subspace** of V if U is also a vector space (with the same addition and scalar multiplication as on V), denoted as $U \leq V$.

The sets $\{0\}$ and V are always subspaces of V . The subspace $\{0\}$ is called the *zero subspace* or *trivial subspace*. Subspaces other than V are called *proper subspaces*.

The following result is useful in determining whether a given subset of V is a subspace of V .

Lemma 4.10 (Subspace test). *Suppose $U \subset V$. Then $U \leq V$ if and only if U satisfies the following conditions:*

- (i) $0 \in U$; (additive identity)
- (ii) $u + w \in U$ for all $u, w \in U$; (closed under addition)
- (iii) $\lambda u \in U$ for all $\lambda \in \mathbf{F}, u \in U$. (closed under scalar multiplication)

Proof.

\Rightarrow If $U \leq V$, then U satisfies the three conditions above by the definition of vector space.

\Leftarrow Suppose U satisfies the three conditions above. (i) ensures that the additive identity of V is in U . (ii) ensures that addition makes sense on U . (iii) ensures that scalar multiplication makes sense on U .

If $u \in U$, then $-u = (-1)u \in U$ by (iii). Hence every element of U has an additive inverse in U .

The other parts of the definition of a vector space, such as associativity and commutativity, are automatically satisfied for U because they hold on the larger space V . Thus U is a vector space and hence is a subspace of V . □

Proposition 4.11. *Suppose $U \leq V$. Then*

- (i) U is a vector space over \mathbf{F} . In fact, the only subsets of V that are vector spaces over \mathbf{F} are the subspaces of V ;
- (ii) if $W \leq U$, then $W \leq V$ (“a subspace of a subspace is a subspace”).

Proof.

- (i) We first check that we have legitimate operations. Since U is closed under addition, the operation $+$ restricted to U gives a map $U \times U \rightarrow U$. Likewise since U is closed under scalar multiplication, that operation restricted to U gives a map $\mathbf{F} \times U \rightarrow U$.

We now check that U satisfies the vector space axioms.

- (i) Commutativity and associativity of addition are inherited from V .
- (ii) There is an additive identity (by the subspace test).

- (iii) There are additive inverses: if $u \in U$ then multiplying by $-1 \in \mathbf{F}$ and shows that $-u = (-1)u \in U$.
- (iv) The remaining four properties are all inherited from V . That is, they apply to general vectors of V and vectors in U are vectors in V .

(ii) This is immediate from the definition of a subspace.

□

Definition 4.12 (Sum of subsets). Suppose $U_1, \dots, U_n \subset V$. The sum of U_1, \dots, U_n is the set of all possible sums of elements of U_1, \dots, U_n :

$$U_1 + \dots + U_n := \{u_1 + \dots + u_n \mid u_i \in U_i\}.$$

Example 4.13. Suppose that $U = \{(x, 0, 0) \in \mathbf{F}^3 \mid x \in \mathbf{F}\}$ and $W = \{(0, y, 0) \in \mathbf{F}^3 \mid y \in \mathbf{F}\}$. Then

$$U + W = \{(x, y, 0) \mid x, y \in \mathbf{F}\}.$$

Suppose that $U = \{(x, x, y, y) \in \mathbf{F}^4 \mid x, y \in \mathbf{F}\}$ and $W = \{(x, x, x, y) \in \mathbf{F}^4 \mid x, y \in \mathbf{F}\}$. Then

$$U + W = \{(x, x, y, z) \in \mathbf{F}^4 \mid x, y, z \in \mathbf{F}\}.$$

The next result states that the sum of subspaces is a subspace, and is in fact the smallest subspace containing all the summands.

Proposition 4.14. Suppose $U_1, \dots, U_n \leq V$. Then $U_1 + \dots + U_n$ is the smallest subspace of V containing U_1, \dots, U_n .

Proof. It is easy to see that $\mathbf{0} \in U_1 + \dots + U_n$ and that $U_1 + \dots + U_n$ is closed under addition and scalar multiplication. Hence by the subspace test, $U_1 + \dots + U_n \leq V$.

Let M be the smallest subspace of V containing U_1, \dots, U_n . We want to show that $U_1 + \dots + U_n = M$. To do so, we show double inclusion: $U_1 + \dots + U_n \subset M$ and $M \subset U_1 + \dots + U_n$.

(i) For all $u_i \in U_i$ ($1 \leq i \leq n$),

$$u_i = \mathbf{0} + \dots + \mathbf{0} + u_i + \mathbf{0} + \dots + \mathbf{0} \in U_1 + \dots + U_n,$$

where all except one of the u 's are $\mathbf{0}$. Thus $U_i \subset U_1 + \dots + U_n$ for $1 \leq i \leq n$. Hence $M \subset U_1 + \dots + U_n$.

(ii) Conversely, every subspace of V containing U_1, \dots, U_n contains $U_1 + \dots + U_n$ (because subspaces must contain all finite sums of their elements). Hence $U_1 + \dots + U_n \subset M$.

□

Definition 4.15 (Direct sum). Suppose $U_1, \dots, U_n \leq V$. If each element of $U_1 + \dots + U_n$ can be written in only one way as a sum $u_1 + \dots + u_n$, $u_i \in U_i$, then $U_1 + \dots + U_n$ is called a **direct sum**. In this case, we denote the sum as

$$U_1 \oplus \dots \oplus U_n.$$

Example 4.16. Suppose that $U = \{(x, y, 0) \in \mathbf{F}^3 \mid x, y \in \mathbf{F}\}$ and $W = \{(0, 0, z) \in \mathbf{F}^3 \mid z \in \mathbf{F}\}$. Then $\mathbf{F}^3 = U \oplus W$.

Suppose U_i is the subspace of \mathbf{F}^n of those vectors whose coordinates are all 0 except for the i -th coordinate; that is, $U_i = \{(0, \dots, 0, x, 0, \dots, 0) \in \mathbf{F}^n \mid x \in \mathbf{F}\}$. Then $\mathbf{F}^n = U_1 \oplus \dots \oplus U_n$.

Lemma 4.17 (Condition for direct sum). *Suppose $V_1, \dots, V_n \leq V$, let $W = V_1 + \dots + V_n$. Then the following are equivalent:*

(i) *Any element in W can be uniquely expressed as the sum of vectors in V_1, \dots, V_n .*

(ii) *If $v_i \in V_i$ satisfies $v_1 + \dots + v_n = \mathbf{0}$, then $v_1 = \dots = v_n = \mathbf{0}$.*

(iii) *For $k = 2, \dots, n$, $(V_1 + \dots + V_{k-1}) \cap V_k = \{\mathbf{0}\}$.*

Proof.

(i) \iff (ii) First suppose W is a direct sum. Then by the definition of direct sum, the only way to write $\mathbf{0}$ as a sum $u_1 + \dots + u_n$ is by taking $u_i = \mathbf{0}$.

Now suppose that the only way to write $\mathbf{0}$ as a sum $v_1 + \dots + v_n$ by taking $v_1 = \dots = v_n = \mathbf{0}$. For $v \in V_1 + \dots + V_n$, suppose that there is more than one way to represent v :

$$\begin{aligned} v &= v_1 + \dots + v_n \\ v &= v'_1 + \dots + v'_n \end{aligned}$$

for some $v_i, v'_i \in V_i$. Subtracting the above two equations gives

$$\mathbf{0} = (v_1 - v'_1) + \dots + (v_n - v'_n).$$

Since $v_i - v'_i \in V_i$, we have $v_i - v'_i = \mathbf{0}$ so $v_i = v'_i$. Hence there is only one unique way to represent $v_1 + \dots + v_n$, thus W is a direct sum.

(ii) \iff (iii) First suppose if $v_i \in V_i$ satisfies $v_1 + \dots + v_n = \mathbf{0}$, then $v_1 = \dots = v_n = \mathbf{0}$. Let $v_k \in (V_1 + \dots + V_{k-1}) \cap V_k$. Then $v_k = v_1 + \dots + v_{k-1}$ where $v_i \in V_i$ ($1 \leq i \leq k-1$). Thus

$$\begin{aligned} v_1 + \dots + v_{k-1} - v_k &= \mathbf{0} \\ v_1 + \dots + v_{k-1} + (-v_k) + \mathbf{0} + \dots + \mathbf{0} &= \mathbf{0} \end{aligned}$$

by taking $v_{k+1} = \dots = v_n = \mathbf{0}$. Then $v_1 = \dots = v_k = \mathbf{0}$.

Now suppose that for $k = 2, \dots, n$, $(V_1 + \dots + V_{k-1}) \cap V_k = \{\mathbf{0}\}$.

$$\begin{aligned} v_1 + \dots + v_n &= \mathbf{0} \\ v_1 + \dots + v_{n-1} &= -v_n \end{aligned}$$

where $v_1 + \dots + v_{n-1} \in V_1 + \dots + V_{n-1}$, $-v_n \in V_n$. Thus

$$v_1 + \dots + v_{n-1} = -v_n \in (V_1 + \dots + V_{n-1}) \cap V_n = \{\mathbf{0}\}$$

so $v_1 + \dots + v_{n-1} = \mathbf{0}$, $v_n = \mathbf{0}$. Induction on n gives $v_1 = \dots = v_{n-1} = v_n = \mathbf{0}$. \square

Proposition 4.18. *Suppose $U, W \leq V$. Then $U + W$ is a direct sum if and only if $U \cap W = \{\mathbf{0}\}$.*

Proof.

\Rightarrow Suppose that $U + W$ is a direct sum. If $v \in U \cap W$, then $\mathbf{0} = v + (-v)$, where $v \in U$, $-v \in W$. By the unique representation of $\mathbf{0}$ as the sum of a vector in U and a vector in W , we have $v = \mathbf{0}$. Thus $U \cap W = \{\mathbf{0}\}$.

\Leftarrow Suppose $U \cap W = \{\mathbf{0}\}$. Suppose $u \in U$, $w \in W$, and $0 = u + w$. $u = -w \in W$, thus $u \in U \cap W$, so $u = w = \mathbf{0}$. By 4.17, $U + W$ is a direct sum. \square

§4.3 Span and Linear Independence

Definition 4.19 (Linear combination). v is a **linear combination** of vectors $v_1, \dots, v_n \in V$ if there exists $a_1, \dots, a_n \in \mathbf{F}$ such that

$$v = a_1v_1 + \dots + a_nv_n.$$

Definition 4.20 (Span). The **span** of $\{v_1, \dots, v_n\}$ is the set of all linear combinations of v_1, \dots, v_n :

$$\text{span}(v_1, \dots, v_n) := \{a_1v_1 + \dots + a_nv_n \mid a_i \in \mathbf{F}\}.$$

The span of the empty set $\{\}$ is defined to be $\{\mathbf{0}\}$.

We say that v_1, \dots, v_n *spans* V if $\text{span}(v_1, \dots, v_n) = V$.

If $S \subset V$ is such that $\text{span}(S) = V$, we say S *spans* V , and S is a *spanning set* for V :

$$\text{span}(S) := \{a_1v_1 + \dots + a_nv_n \mid v_i \in S, a_i \in \mathbf{F}\}.$$

Proposition 4.21. $\text{span}(v_1, \dots, v_n)$ in V is the smallest subspace of V containing v_1, \dots, v_n .

Proof. First we show that $\text{span}(v_1, \dots, v_n) \leq V$, using the subspace test.

- (i) $\mathbf{0} = 0v_1 + \dots + 0v_n \in \text{span}(v_1, \dots, v_n)$
- (ii) $(a_1v_1 + \dots + a_nv_n) + (c_1v_1 + \dots + c_nv_n) = (a_1 + c_1)v_1 + \dots + (a_n + c_n)v_n \in \text{span}(v_1, \dots, v_n)$, so $\text{span}(v_1, \dots, v_n)$ is closed under addition.
- (iii) $\lambda(a_1v_1 + \dots + a_nv_n) = (\lambda a_1)v_1 + \dots + (\lambda a_n)v_n \in \text{span}(v_1, \dots, v_n)$, so $\text{span}(v_1, \dots, v_n)$ is closed under scalar multiplication.

Let M be the smallest vector subspace of V containing v_1, \dots, v_n . We claim that $M = \text{span}(v_1, \dots, v_n)$. To show this, we show that (i) $M \subset \text{span}(v_1, \dots, v_n)$ and (ii) $M \supset \text{span}(v_1, \dots, v_n)$.

- (i) Each v_i is a linear combination of v_1, \dots, v_n , as

$$v_i = 0 \cdot v_1 + \dots + 0 \cdot v_{i-1} + 1 \cdot v_i + 0 \cdot v_{i+1} + \dots + 0 \cdot v_n,$$

so by the definition of the span as the collection of all linear combinations of v_1, \dots, v_n , we have that $v_i \in \text{span}(v_1, \dots, v_n)$. But M is the smallest vector subspace containing v_1, \dots, v_n , so

$$M \subset \text{span}(v_1, \dots, v_n).$$

- (ii) Since $v_i \in M$ ($1 \leq i \leq n$) and M is a vector subspace (closed under addition and scalar multiplication), it follows that

$$a_1v_1 + \dots + a_nv_n \in M$$

for all $a_i \in \mathbf{F}$ (i.e. M contains all linear combinations of v_1, \dots, v_n). So

$$\text{span}(v_1, \dots, v_n) \subset M.$$

□

Definition 4.22 (Finite-dimensional). V is *finite-dimensional* if there exists some list of vector $\{v_1, \dots, v_n\}$ that spans V ; otherwise, it is *infinite-dimensional*.

Remark. Recall that by definition every list of vectors has finite length.

Remark. From this definition, infinite-dimensionality is the negation of finite-dimensionality (i.e. *not* finite-dimensional). Hence to prove that a vector space is infinite-dimensional, we prove by contradiction; that is, first assume that the vector space is finite-dimensional, then try to come to a contradiction.

Exercise 4.1. For positive integer n , \mathbf{F}^n is finite-dimensional.

Proof. Suppose $(x_1, x_2, \dots, x_n) \in \mathbf{F}^n$, then

$$(x_1, x_2, \dots, x_n) = x_1(1, 0, \dots, 0) + x_2(0, 1, \dots, 0) + \dots + x_n(0, 0, \dots, 1)$$

so

$$(x_1, \dots, x_n) \in \text{span}((1, 0, \dots, 0), (0, 1, \dots, 0), \dots, (0, \dots, 0, 1)).$$

The vectors $(1, 0, \dots, 0), (0, 1, \dots, 0), \dots, (0, \dots, 0, 1)$ spans \mathbf{F}^n , so \mathbf{F}^n is finite-dimensional. □

Definition 4.23 (Linear independence). A list of vectors v_1, \dots, v_n is *linearly independent* in V if the only choice of $a_1, \dots, a_n \in \mathbf{F}$ that makes

$$a_1v_1 + \dots + a_nv_n = \mathbf{0}$$

is $a_1 = \dots = a_n = 0$; otherwise, it is *linearly dependent*.

We say that $S \subset V$ is linearly independent if every finite subset of S is linearly independent.

Lemma 4.24 (Compare coefficients). *Let v_1, \dots, v_n be linearly independent in V . Then*

$$a_1v_1 + \dots + a_nv_n = b_1v_1 + \dots + b_nv_n$$

if and only if $a_i = b_i$ ($1 \leq i \leq n$).

Proof. Exercise. □

The following result will often be useful; it states that given a linearly dependent set of vectors, one of the vectors is in the span of the previous ones; furthermore we can throw out that vector without changing the span of the original set.

Lemma 4.25 (Linear dependence lemma). *Suppose v_1, \dots, v_n are linearly dependent in V . Then there exists v_k such that the following hold:*

(i) $v_k \in \text{span}(v_1, \dots, v_{k-1})$

(ii) $\text{span}(v_1, \dots, v_{k-1}, v_{k+1}, \dots, v_n) = \text{span}(v_1, \dots, v_n)$

Proof.

(i) Since v_1, \dots, v_n are linearly dependent, there exists $a_1, \dots, a_n \in \mathbf{F}$, not all 0, such that

$$a_1v_1 + \dots + a_nv_n = 0.$$

Take $k = \max\{1, \dots, n\}$ such that $a_k \neq 0$. Then

$$v_k = -\frac{a_1}{a_k}v_1 - \dots - \frac{a_{k-1}}{a_k}v_{k-1},$$

which means that v_k can be written as a linear combination of v_1, \dots, v_{k-1} , so $v_k \in \text{span}(v_1, \dots, v_{k-1})$ by definition of span.

(ii) Now suppose k is such that $v_k \in \text{span}(v_1, \dots, v_{k-1})$. Then there exists $b_1, \dots, b_{k-1} \in \mathbf{F}$ be such that

$$v_k = b_1v_1 + \dots + b_{k-1}v_{k-1}. \quad (1)$$

Suppose $u \in \text{span}(v_1, \dots, v_n)$. Then there exists $c_1, \dots, c_n \in \mathbf{F}$ such that

$$u = c_1v_1 + \dots + c_nv_n. \quad (2)$$

In (2), we can replace v_k with the RHS of (1), which gives

$$\begin{aligned} u &= c_1v_1 + \dots + c_{k-1}v_{k-1} + c_kv_k + c_{k+1}v_{k+1} + \dots + c_nv_n \\ &= c_1v_1 + \dots + c_{k-1}v_{k-1} + c_k(b_1v_1 + \dots + b_{k-1}v_{k-1}) + c_{k+1}v_{k+1} + \dots + c_nv_n \\ &= c_1v_1 + \dots + c_{k-1}v_{k-1} + c_kb_1v_1 + \dots + c_kb_{k-1}v_{k-1} + c_{k+1}v_{k+1} + \dots + c_nv_n \\ &= (c_1 + bc_k)v_1 + \dots + (c_{k-1} + b_{k-1}c_k)v_{k-1} + c_{k+1}v_{k+1} + \dots + c_nv_n. \end{aligned}$$

Thus $u \in \text{span}(v_1, \dots, v_{k-1}, v_{k+1}, \dots, v_n)$. This shows that removing v_k from v_1, \dots, v_n does not change the span of the list.

□

The following result says that no linearly independent set in V is longer than a spanning set in V .

Proposition 4.26. *In a finite-dimensional vector space, the length of every linearly independent set of vectors is less than or equal to the length of every spanning set of vectors.*

Proof. Suppose $A = \{u_1, \dots, u_m\}$ is linearly independent in V , $B = \{w_1, \dots, w_n\}$ spans V . We want to prove that $m \leq n$.

Since B spans V , if we add any other vector from V to the list B , we will get a linearly dependent list, since this newly added vector can, by the definition of a span, be expressed as a linear combination of the vectors in B . In particular, if we add $u_1 \in A$ to B , then the new list

$$\{u_1, w_1, \dots, w_n\}$$

is linearly dependent. By the linear independence lemma, we can remove one of the w_i 's from B , so that the remaining list of n vectors still spans V . For the sake of argument, let's say we remove w_n (we can always order

the w_i 's in the list so that the element we remove is at the end). Then we are left with the revised list

$$B_1 = \{u_1, w_1, \dots, w_{n-1}\}.$$

We can repeat this process m times, each time adding the next element u_i from list A and removing the last w_i . Because of the linear dependence lemma, we know that there must always be a w_i that can be removed each time we add a u_i , so there must be at least as many w_i 's as u_i 's. In other words, $m \leq n$ which is what we wanted to prove. \square

Remark. We can use this result to show, without any computations, that certain lists are not linearly independent and that certain lists do not span a given vector space.

Our intuition suggests that every subspace of a finite-dimensional vector space should also be finite-dimensional. We now prove that this intuition is correct.

Proposition 4.27. *Every subspace of a finite-dimensional vector space is finite-dimensional.*

Proof. Suppose V is finite-dimensional, $U \leq V$. To show that U is finite-dimensional, we need to find a spanning set of vectors in U . We prove by construction of this spanning set.

Step 1 If $U = \{\mathbf{0}\}$, then U is finite-dimensional and we are done. Otherwise, choose $v_1 \in U$, $v_1 \neq \mathbf{0}$ and add it to our list of vectors.

Step k Our list so far is $\{v_1, \dots, v_{k-1}\}$. If $U = \text{span}(v_1, \dots, v_{k-1})$, then U is finite-dimensional and we are done. Otherwise, choose $v_k \in U$ such that $v_k \notin \text{span}(v_1, \dots, v_{k-1})$ and add it to our list.

After each step, we have constructed a list of vectors such that no vector in this list is in the span of the previous vectors; by the linear dependence lemma, our constructed list is a linearly independent set.

By 4.26, this linearly independent set cannot be longer than any spanning set of V . Thus the process must terminate after a finite number of steps, and we have constructed a spanning set of U . Hence U is finite-dimensional. \square

§4.4 Bases

Definition 4.28 (Basis). $B = \{v_1, \dots, v_n\}$ is a *basis* of V if

- (i) B is linearly independent in V ;
- (ii) B is a spanning set of V .

Example 4.29 (Standard basis). Let $e_i = (0, \dots, 0, 1, 0, \dots, 0)$ where the i -th coordinate is 1. $\{e_1, \dots, e_n\}$ is a basis of \mathbf{F}^n , known as the *standard basis* of \mathbf{F}^n .

Lemma 4.30 (Criterion for basis). Let $B = \{v_1, \dots, v_n\}$ be a set of vectors in V . Then B is a basis of V if and only if every $v \in V$ can uniquely expressed as a linear combination of v_1, \dots, v_n .

Proof.

\Rightarrow Let $v \in V$. Since B is a basis of V , there exists $a_1, \dots, a_n \in \mathbf{F}$ such that

$$v = a_1v_1 + \dots + a_nv_n. \quad (1)$$

To show that the representation is unique, suppose that $c_1, \dots, c_n \in \mathbf{F}$ also satisfy

$$v = c_1v_1 + \dots + c_nv_n. \quad (2)$$

Subtracting (2) from (1) gives

$$\mathbf{0} = (a_1 - c_1)v_1 + \dots + (a_n - c_n)v_n.$$

Since v_1, \dots, v_n are linearly independent, we have $a_i - c_i = 0$, or $a_i = c_i$ for all i . Thus the representation of v as a linear combination of v_1, \dots, v_n is unique.

\Leftarrow Suppose that every $v \in V$ can be uniquely expressed as a linear combination of v_1, \dots, v_n . This implies that B spans V .

To show that B is linearly independent, suppose that $a_1, \dots, a_n \in \mathbf{F}$ are such that

$$a_1v_1 + \dots + a_nv_n = \mathbf{0}.$$

Since $\mathbf{0}$ can be uniquely expressed as a linear combination of v_1, \dots, v_n , we have $a_1 = \dots = a_n = 0$, thus B is linearly independent. Since B is linearly independent and spans V , B is a basis of V . \square

A spanning set in a vector space may not be a basis because it is not linearly independent. The next result says that given any spanning set, some (possibly none) of the vectors in it can be discarded so that the remaining list is linearly independent and still spans the vector space.

Lemma 4.31. Every spanning set in a vector space can be reduced to a basis of the vector space.

Proof. Suppose $B = \{v_1, \dots, v_n\}$ spans V . We want to remove some vectors from B so that the remaining vectors form a basis of V . We do this through the multistep process described below.

Step 1 If $v_1 = \mathbf{0}$, delete v_1 from B . If $v_1 \neq \mathbf{0}$, leave B unchanged.

Step k If $v_k \in \text{span}(v_1, \dots, v_{k-1})$, delete v_k from B . If $v_k \notin \text{span}(v_1, \dots, v_{k-1})$, leave B unchanged.

Stop the process after step n , getting a list B . Since we only delete vectors from B that are in the span of the previous vectors, by the linear dependence lemma, the list B still spans V .

The process ensures that no vector in B is in the span of the previous ones. By the linear dependence lemma, B is linearly independent.

Since B is linearly independent and spans V , B is a basis of V . \square

Corollary 4.32. *Every finite-dimensional vector space has a basis.*

Proof. We prove by construction. Suppose V is finite-dimensional. By definition, there exists a spanning set of vectors in V . By 4.31, the spanning set can be reduced to a basis. \square

Now we show that given any linearly independent set, we can adjoin some additional vectors so that the extended list is still linearly independent but also spans the space.

Lemma 4.33. *Every linearly independent set of vectors in a finite-dimensional vector space can be extended to a basis of the vector space.*

Proof. Suppose u_1, \dots, u_m are linearly independent in V , w_1, \dots, w_n span V . Then the list

$$\{u_1, \dots, u_m, w_1, \dots, w_n\}$$

spans V . By 4.31, we can reduce this list to a basis of V consisting u_1, \dots, u_m (since u_1, \dots, u_m are linearly independent, $u_i \notin \text{span}(u_1, \dots, u_{i-1})$ for all i , so none of the u_i 's are deleted in the process), and some of the w_i 's. \square

We now show that every subspace of a finite-dimensional vector space can be paired with another subspace to form a direct sum of the whole space.

Corollary 4.34. *Suppose V is finite-dimensional, $U \leq V$. Then there exists $W \leq V$ such that $V = U \oplus W$.*

Proof. Since V is finite-dimensional and $U \leq V$, by 4.27, U is finite-dimensional; by 4.32, U has a basis B . Let $B = \{u_1, \dots, u_n\}$.

Since B is linearly independent, by 4.33, B can be extended to a basis of V , say

$$\{u_1, \dots, u_n, w_1, \dots, w_n\}.$$

Claim. $W = \text{span}(w_1, \dots, w_n)$.

We need to show that $V = U \oplus W$; by 4.17, we need to show (i) $V = U + W$, and (ii) $U \cap W = \{0\}$.

(i) Let $v \in V$. Since $\{u_1, \dots, u_n, w_1, \dots, w_n\}$ spans V , there exists $a_1, \dots, a_n, b_1, \dots, b_n \in \mathbf{F}$ such that

$$v = a_1 u_1 + \dots + a_n u_n + b_1 w_1 + \dots + b_n w_n.$$

Take $u = a_1u_1 + \cdots + a_nu_n \in U$, $w = b_1w_1 + \cdots + b_nw_n \in W$. Then $v = u + w \in U + W$, so $V = U + W$.

(ii) Let $v \in U \cap W$. Since $v \in U$, v can be written as a linear combination of u_1, \dots, u_n :

$$v = a_1u_1 + \cdots + a_nu_n. \quad (1)$$

Since $v \in W$, v can be written as a linear combination of w_1, \dots, w_n :

$$v = b_1w_1 + \cdots + b_nw_n. \quad (2)$$

Subtracting (2) from (1) gives

$$\mathbf{0} = a_1u_1 + \cdots + a_nu_n - b_1w_1 - \cdots - b_nw_n.$$

Since $u_1, \dots, u_n, w_1, \dots, w_n$ are linearly independent, we have $a_i = b_i = 0$ for all i ($1 \leq i \leq n$). Thus $v = \mathbf{0}$, so $U \cap W = \{\mathbf{0}\}$.

□

§4.5 Dimension

Lemma 4.35. *Any two bases of a finite-dimensional vector space have the same length.*

Proof. Suppose V is finite-dimensional, let B_1 and B_2 be two bases of V . By definition, B_1 is linearly independent in V , and B_2 spans V , so by 4.26, $|B_1| \leq |B_2|$.

Similarly, by definition, B_2 is linearly independent in V and B_1 spans V , so $|B_2| \leq |B_1|$.

Since $|B_1| \leq |B_2|$ and $|B_2| \leq |B_1|$, we have $|B_1| = |B_2|$, as desired. \square

Since any two bases of a finite-dimensional vector space have the same length, we can formally define the dimension of such spaces.

Definition 4.36 (Dimension). The *dimension* of V is the length of any basis of V , denoted by $\dim V$.

Proposition 4.37. *Suppose V is finite-dimensional, $U \leq V$. Then $\dim U \leq \dim V$.*

Proof. Since V is finite-dimensional and $U \leq V$, U is finite-dimensional. Let B_U be a basis of U , and B_V be a basis of V .

By definition, B_U is linearly independent in V , and B_V spans V . By 4.26, $|B_U| \leq |B_V|$, so

$$\dim U = |B_U| \leq |B_V| = \dim V,$$

since $|B_U| = \dim U$ and $|B_V| = \dim V$ by definition. \square

To check that a list of vectors is a basis, we must show that it is linearly independent and that it spans the vector space. The next result shows that if the list in question has the right length, then we only need to check that it satisfies one of the two required properties.

Proposition 4.38. *Suppose V is finite-dimensional. Then*

- (i) *every linearly independent set of vectors in V with length $\dim V$ is a basis of V ;*
- (ii) *every spanning set of vectors in V with length $\dim V$ is a basis of V .*

Proof.

- (i) Suppose $\dim V = n$, $\{v_1, \dots, v_n\}$ is linearly independent in V . By 4.33, $\{v_1, \dots, v_n\}$ can be extended to a basis of V . However, every basis of V has length n (by definition of dimension), which means that no elements are adjoined to $\{v_1, \dots, v_n\}$. Hence $\{v_1, \dots, v_n\}$ is a basis of V , as desired.
- (ii) Suppose $\dim V = n$, $\{v_1, \dots, v_n\}$ spans V . By 4.31, $\{v_1, \dots, v_n\}$ can be reduced to a basis of V . However, every basis of V has length n , which means that no elements are deleted from $\{v_1, \dots, v_n\}$. Hence $\{v_1, \dots, v_n\}$ is a basis of V , as desired.

\square

Corollary 4.39. *Suppose V is finite-dimensional, $U \leq V$. If $\dim U = \dim V$, then $U = V$.*

Proof. Let $\dim U = \dim V = n$, let $\{u_1, \dots, u_n\}$ be a basis of U . Then $\{u_1, \dots, u_n\}$ is linearly independent in V (because it is a basis of U) of length $\dim V$. From 4.38, $\{u_1, \dots, u_n\}$ is a basis of V . In particular every vector in V is a linear combination of u_1, \dots, u_n . Thus $U = V$. \square

Lemma 4.40 (Dimension of sum). *Suppose V is finite-dimensional, $U_1, U_2 \leq V$. Then*

$$\dim(U_1 + U_2) = \dim U_1 + \dim U_2 - \dim(U_1 \cap U_2).$$

Proof. Let $\{u_1, \dots, u_m\}$ be a basis of $U_1 \cap U_2$; thus $\dim(U_1 \cap U_2) = m$. Since $\{u_1, \dots, u_m\}$ is a basis of $U_1 \cap U_2$, it is linearly independent in U_1 . By 4.33, $\{u_1, \dots, u_m\}$ can be extended to a basis $\{u_1, \dots, u_m, v_1, \dots, v_j\}$ of U_1 ; thus $\dim U_1 = m + j$. Similarly, extend $\{u_1, \dots, u_m\}$ to a basis $\{u_1, \dots, u_m, v_1, \dots, v_k\}$ of U_2 ; thus $\dim U_2 = m + k$.

We will show that

$$\{u_1, \dots, u_m, v_1, \dots, v_j, w_1, \dots, w_k\}$$

is a basis of $U_1 + U_2$. This will complete the proof because then we will have

$$\begin{aligned} \dim(U_1 + U_2) &= m + j + k \\ &= (m + j) + (m + k) - m \\ &= \dim U_1 + \dim U_2 - \dim(U_1 \cap U_2). \end{aligned}$$

We just need to show that $\{u_1, \dots, u_m, v_1, \dots, v_j, w_1, \dots, w_k\}$ is linearly independent. To prove this, suppose

$$a_1 u_1 + \dots + a_m u_m + b_1 v_1 + \dots + b_j v_j + c_1 w_1 + \dots + c_k w_k = \mathbf{0}, \quad (1)$$

where $a_i, b_i, c_i \in \mathbf{F}$. We need to show that $a_i = b_i = c_i = 0$ for all i . (1) can be rewritten as

$$c_1 w_1 + \dots + c_k w_k = -a_1 u_1 - \dots - a_m u_m - b_1 v_1 - \dots - b_j v_j,$$

which shows that $c_1 w_1 + \dots + c_k w_k \in U_1$. But actually all the w_i 's are in U_2 , so $c_1 w_1 + \dots + c_k w_k \in U_2$, thus $c_1 w_1 + \dots + c_k w_k \in U_1 \cap U_2$. Then we can write

$$c_1 w_1 + \dots + c_k w_k = d_1 u_1 + \dots + d_m u_m$$

for some $d_i \in \mathbf{F}$. But $u_1, \dots, u_m, w_1, \dots, w_k$ are linearly independent, so $c_i = d_i = 0$ for all i . Thus our original equation (1) becomes

$$a_1 u_1 + \dots + a_m u_m + b_1 v_1 + \dots + b_j v_j = \mathbf{0}.$$

Since $u_1, \dots, u_m, v_1, \dots, v_j$ are linearly independent, $a_i = b_i = 0$ for all i , as desired. \square

Exercises

Exercise 4.2 ([Ax124] 1C Q12). Suppose W is a vector space over \mathbf{F} , V_1 and V_2 are subspaces of W . Show that $V_1 \cup V_2$ is a vector space over \mathbf{F} if and only if $V_1 \subset V_2$ or $V_2 \subset V_1$.

Solution. The backward direction is trivial. We focus on proving the forward direction.

Supppse otherwise, then $V_1 \setminus V_2 \neq \emptyset$ and $V_2 \setminus V_1 \neq \emptyset$. Pick $v_1 \in V_1 \setminus V_2$ and $v_2 \in V_2 \setminus V_1$. Then

$$\begin{aligned} v_1, v_2 \in V_1 \cup V_2 &\implies v_1 + v_2 \in V_1 \cup V_2 \\ &\implies v_2, v_1 + v_2 \in V_2 \\ &\implies v_1 = (v_1 + v_2) - v_2 \in V_2 \end{aligned}$$

which is a contradiction. □

Exercise 4.3 ([Ax124] 1C Q13). Suppose W is a vector space over \mathbf{F} , V_1, V_2, V_3 are subspaces of W . Then $V_1 \cup V_2 \cup V_3$ is a vector space over \mathbf{F} if and only if one of the V_i contains the other two.

Solution. We prove the forward direction. Suppose otherwise, then $v_1 \in V_1 \setminus (V_2 + V_3)$, $v_2 \in V_2 \setminus (V_1 + V_3)$, $v_3 \in V_3 \setminus (V_1 + V_2)$. Consider

$$\{v_1 + v_2 + v_3, v_1 + v_2 + 2v_3, v_1 + 2v_2 + v_3, v_1 + 2v_2 + 2v_3\} \subset V_1 \cup V_2 \cup V_3$$

Then

$$\begin{aligned} (v_1 + v_2 + 2v_3) - (v_1 + v_2 + v_3) &= v_3 \notin V_1 + V_2 \\ \implies v_1 + v_2 + v_3 &\notin V_1 + V_2 \quad \text{or} \quad v_1 + v_2 + 2v_3 \notin V_1 + V_2 \\ \implies v_1 + v_2 + v_3 &\in V_3 \quad \text{or} \quad v_1 + v_2 + 2v_3 \in V_3 \\ \implies v_1 + v_2 &\in V_3 \end{aligned}$$

Similarly,

$$\begin{aligned} (v_1 + 2v_2 + 2v_3) - (v_1 + 2v_2 + v_3) &= v_3 \notin V_1 + V_2 \\ \implies v_1 + 2v_2 + v_3 &\notin V_1 + V_2 \quad \text{or} \quad v_1 + 2v_2 + 2v_3 \notin V_1 + V_2 \\ \implies v_1 + 2v_2 + v_3 &\in V_3 \quad \text{or} \quad v_1 + 2v_2 + 2v_3 \in V_3 \\ \implies v_1 + 2v_2 &\in V_3 \end{aligned}$$

This implies $(v_1 + 2v_2) - (v_1 + v_2) = v_2 \in V_3$, a contradiction. □

Exercise 4.4 ([Ax124] 2A Q12). Suppose $\{v_1, \dots, v_n\}$ is linearly independent in V , $w \in V$. Prove that if $\{v_1 + w, \dots, v_n + w\}$ is linearly dependent, then $w \in \text{span}(v_1, \dots, v_n)$.

Solution. If $\{v_1 + w, \dots, v_n + w\}$ is linearly dependent, then there exists $a_1, \dots, a_n \in \mathbf{F}$, not all zero, such that

$$a_1(v_1 + w) + \dots + a_n(v_n + w) = 0,$$

or

$$a_1v_1 + \dots + a_nv_n = -(a_1 + \dots + a_n)w.$$

Suppose otherwise, that $a_1 + \cdots + a_n = 0$. Then

$$a_1v_1 + \cdots + a_nv_n = \mathbf{0},$$

but the linear independence of $\{v_1, \dots, v_n\}$ implies that $a_1 = \cdots = a_n = 0$, which is a contradiction. Hence we must have $a_1 + \cdots + a_n \neq 0$, so we can write

$$w = -\frac{a_1}{a_1 + \cdots + a_n}v_1 - \cdots - \frac{a_n}{a_1 + \cdots + a_n}v_n,$$

which is a linear combination of v_1, \dots, v_n . Thus by definition of span, $w \in \text{span}(v_1, \dots, v_n)$. \square

Exercise 4.5 ([Ax124] 2A Q14). Suppose $\{v_1, \dots, v_n\} \subset V$. Let

$$w_i = v_1 + \cdots + v_i \quad (i = 1, \dots, n)$$

Show that $\{v_1, \dots, v_n\}$ is linearly independent if and only if $\{w_1, \dots, w_n\}$ is linearly independent.

Solution. Write

$$\begin{aligned} v_1 &= w_1 \\ v_2 &= w_2 - w_1 \\ v_3 &= w_3 - w_2 \\ &\vdots \\ v_n &= w_n - w_{n-1}. \end{aligned}$$

\Rightarrow

$$a_1w_1 + \cdots + a_nw_n = \mathbf{0}$$

for some $a_i \in \mathbf{F}$. Expressing w_i 's as v_i 's,

$$a_1v_1 + a_2(v_1 + v_2) + \cdots + a_n(v_1 + \cdots + v_n) = \mathbf{0},$$

or

$$(a_1 + \cdots + a_n)v_1 + (a_2 + \cdots + a_n)v_2 + \cdots + a_nv_n = \mathbf{0}.$$

Since v_1, \dots, v_n are linearly independent,

$$\begin{aligned} a_1 + a_2 + \cdots + a_n &= 0 \\ a_2 + \cdots + a_n &= 0 \\ &\vdots \\ a_n &= 0 \end{aligned}$$

on solving simultaneously gives $a_1 = \cdots = a_n = 0$.

\Leftarrow Similar to the above. \square

Exercise 4.6 ([Ax124] 2A Q18). Prove that \mathbf{F}^∞ is infinite-dimensional.

Solution. To prove that \mathbf{F}^∞ has no finite spanning sets, we prove by contradiction. Suppose otherwise, that there exists a finite spanning set of \mathbf{F}^∞ , say $\{v_1, \dots, v_n\}$.

Let

$$\begin{aligned} e_1 &= (1, 0, \dots) \\ e_2 &= (0, 1, 0, \dots) \\ e_3 &= (0, 0, 1, 0, \dots) \\ &\vdots \\ e_{n+1} &= (0, \dots, 0, 1, 0, \dots) \end{aligned}$$

where e_i has a 1 at the i -th coordinate, and 0's for the remaining coordinates. Let

$$a_1 e_1 + \dots + a_{n+1} e_{n+1} = \mathbf{0}$$

for some $a_i \in \mathbf{F}$. Then

$$(a_1, a_2, \dots, a_{n+1}, 0, 0, \dots) = \mathbf{0}$$

so $a_1 = a_2 = \dots = a_{n+1} = 0$. Thus $\{e_1, \dots, e_{n+1}\}$ is a linearly independent set, of length $n + 1$. However, $\{v_1, \dots, v_n\}$ is a spanning set of length n . By 4.26, we have reached a contradiction. \square

Exercise 4.7 ([Axl24] 2B Q5). Suppose V is finite-dimensional, $U, W \leq V$ such that $V = U + W$. Prove that V has a basis in $U \cup W$.

Solution. Let $\{v_i\}_{i=1}^n$ denote the basis for V . By definition we have $v_i = u_i + w_i$ for some $u_i \in U, w_i \in W$. Then we have the spanning set of the vector space V $\sum_{i=1}^n a_i(u_i + w_i)$, which can be reduced to a basis by the lemma. \square

Exercise 4.8 ([Axl24] 2B Q7). Suppose $\{v_1, v_2, v_3, v_4\}$ is a basis of V . Prove that

$$\{v_1 + v_2, v_2 + v_3, v_3 + v_4, v_4\}$$

is also a basis of V .

Solution. We know that $\{v_1, v_2, v_3, v_4\}$ is linearly independent and spans V . Then there exist $a_i \in \mathbf{F}$ such that

$$a_1(v_1 + v_2) + a_2(v_2 + v_3) + a_3(v_3 + v_4) + a_4 v_4 = 0 \implies a_1 = a_2 = a_3 = a_4 = 0.$$

Write

$$\begin{aligned} &a_1(v_1 + v_2) + a_2(v_2 + v_3) + a_3(v_3 + v_4) + a_4 v_4 \\ &= a_1 v_1 + (a_1 + a_2)v_2 + (a_2 + a_3)v_3 + (a_3 + a_4)v_4, \end{aligned}$$

this shows the linear independence. To prove spanning, let $v \in V$, then

$$\begin{aligned} v &= a_1 v_1 + a_2 v_2 + a_3 v_3 + a_4 v_4 \\ &= a_1(v_1 + v_2) + (a_2 - a_1)(v_2 + v_3) + (a_3 - a_2)(v_3 + v_4) + (a_4 - a_3)v_4, \end{aligned}$$

which is a linear combination of $v_1 + v_2, v_2 + v_3, v_3 + v_4, v_4$. \square

Exercise 4.9 ([Ax124] 2B Q10). Suppose $U, W \leq V$ such that $V = U \oplus W$. Suppose also that $\{u_1, \dots, u_m\}$ is a basis of U , $\{w_1, \dots, w_n\}$ is a basis of W . Prove that

$$\{u_1, \dots, u_m, w_1, \dots, w_n\}$$

is a basis of V .

Solution. We know that this set is linearly independent (otherwise violating the direct sum assumption) so it suffices to prove the spanning. Let $v \in V$, then

$$v = u + w = \sum_{i=1}^m a_i u_i + \sum_{j=1}^n b_j w_j.$$

□

Exercise 4.10 ([Ax124] 2C Q8).

Exercise 4.11 ([Ax124] 2C Q16).

Exercise 4.12 ([Ax124] 2C Q17). Suppose that $V_1, \dots, V_n \leq V$ are finite-dimensional. Prove that $V_1 + \dots + V_n$ is finite-dimensional, and

$$\dim(V_1 + \dots + V_n) \leq \dim V_1 + \dots + \dim V_n.$$

Solution. We prove by induction on n . The base case is trivial. Assume the statement holds for k . Then for $k + 1$, denoting $V_1 + \dots + V_k = M_k$, we have that

$$\dim(M_k + V_{k+1}) \leq \dim V_1 + \dots + \dim V_{k+1},$$

which is finite.

□

5 Linear Maps

§5.1 Vector Space of Linear Maps

Definition 5.1 (Linear map). A **linear map** from V to W is a function $T: V \rightarrow W$ satisfying the following properties:

- (i) $T(v + w) = Tv + Tw$ for all $v, w \in V$; (additivity)
- (ii) $T(\lambda v) = \lambda T(v)$ for all $\lambda \in \mathbf{F}, v \in V$. (homogeneity)

Notation. The set of linear maps from V to W is denoted by $\mathcal{L}(V, W)$; the set of linear maps on V (from V to V) is denoted by $\mathcal{L}(V)$.

The existence part of the next result means that we can find a linear map that takes on whatever values we wish on the vectors in a basis. The uniqueness part of the next result means that a linear map is completely determined by its values on a basis.

Lemma 5.2 (Linear map lemma). *Suppose $\{v_1, \dots, v_n\}$ is a basis of V , and $w_1, \dots, w_n \in W$. Then there exists a unique linear map $T: V \rightarrow W$ such that*

$$Tv_i = w_i \quad (i = 1, \dots, n)$$

Proof. First we show the existence of a linear map T with the desired property. Define $T: V \rightarrow W$ by

$$T(c_1v_1 + \dots + c_nv_n) = c_1w_1 + \dots + c_nw_n,$$

for some $c_i \in \mathbf{F}$. Since $\{v_1, \dots, v_n\}$ is a basis of V , by 4.30, each $v \in V$ can be uniquely expressed as a linear combination of v_1, \dots, v_n , thus the equation above does indeed define a function $T: V \rightarrow W$. For i ($1 \leq i \leq n$), take $c_i = 1$ and the other c 's equal to 0, then

$$T(0v_1 + \dots + 1v_i + \dots + 0v_n) = 0w_1 + \dots + 1w_i + \dots + 0w_n$$

which shows that $Tv_i = w_i$.

We now show that $T: V \rightarrow W$ is a linear map:

- (i) For $u, v \in V$ with $u = a_1v_1 + \dots + a_nv_n$ and $v = c_1v_1 + \dots + c_nv_n$,

$$\begin{aligned} T(u + v) &= T((a_1 + c_1)v_1 + \dots + (a_n + c_n)v_n) \\ &= (a_1 + c_1)w_1 + \dots + (a_n + c_n)w_n \\ &= (a_1w_1 + \dots + a_nw_n) + (c_1w_1 + \dots + c_nw_n) \\ &= Tu + Tv. \end{aligned}$$

(ii) For $\lambda \in \mathbf{F}$ and $v = c_1v_1 + \cdots + c_nv_n$,

$$\begin{aligned} T(\lambda v) &= T(\lambda c_1v_1 + \cdots + \lambda c_nv_n) \\ &= \lambda c_1w_1 + \cdots + \lambda c_nw_n \\ &= \lambda(c_1w_1 + \cdots + c_nw_n) \\ &= \lambda Tv. \end{aligned}$$

To prove uniqueness, now suppose that $T \in \mathcal{L}(V, W)$ and $Tv_i = w_i$ for $i = 1, \dots, n$. Let $c_i \in \mathbf{F}$. The homogeneity of T implies that $T(c_iv_i) = c_iw_i$. The additivity of T now implies that

$$T(c_1v_1 + \cdots + c_nv_n) = c_1w_1 + \cdots + c_nw_n.$$

Thus T is uniquely determined on $\text{span}\{v_1, \dots, v_n\}$. Since $\{v_1, \dots, v_n\}$ is a basis of V , this implies that T is uniquely determined on V . \square

Proposition 5.3. $\mathcal{L}(V, W)$ is a vector space, with the operations addition and scalar multiplication defined as follows: suppose $S, T \in \mathcal{L}(V, W)$, $\lambda \in \mathbf{F}$,

$$(i) (S + T)(v) = Sv + Tv$$

$$(ii) (\lambda T)(v) = \lambda(Tv)$$

for all $v \in V$.

Proof. Exercise. \square

Definition 5.4 (Product of linear maps). $T \in \mathcal{L}(U, V)$, $S \in \mathcal{L}(V, W)$, then the **product** $ST \in \mathcal{L}(U, W)$ is defined by

$$(ST)(u) = S(Tu) \quad (u \in U)$$

Remark. In other words, ST is just the usual composition $S \circ T$ of two functions.

Remark. ST is defined only when T maps into the domain of S .

Proposition 5.5 (Algebraic properties of products of linear maps).

(i) *Associativity:* $(T_1T_2)T_3 = T_1(T_2T_3)$ for all linear maps T_1, T_2, T_3 such that the products make sense (meaning that T_3 maps into the domain of T_2 , T_2 maps into the domain of T_1)

(ii) *Identity:* $TI = IT = T$ for all $T \in \mathcal{L}(V, W)$ (the first I is the identity map on V , and the second I is the identity map on W)

(iii) *Distributive:* $(S_1 + S_2)T = S_1T + S_2T$ and $S(T_1 + T_2) = ST_1 + ST_2$ for all $T, T_1, T_2 \in \mathcal{L}(U, V)$ and $S, S_1, S_2 \in \mathcal{L}(V, W)$

Proof. Exercise. \square

Proposition 5.6. Suppose $T \in \mathcal{L}(V, W)$. Then $T(\mathbf{0}) = \mathbf{0}$.

Proof. By additivity, we have

$$T(\mathbf{0}) = T(\mathbf{0} + \mathbf{0}) = T(\mathbf{0}) + T(\mathbf{0}).$$

Add the additive inverse of $T(\mathbf{0})$ to each side of the equation to conclude that $T(\mathbf{0}) = \mathbf{0}$.

□

§5.2 Kernel and Image

Definition 5.7 (Kernel). Suppose $T \in \mathcal{L}(V, W)$. The *kernel* of T is the subset of V consisting of those vectors that T maps to $\mathbf{0}$:

$$\ker T := \{v \in V \mid Tv = \mathbf{0}\} \subset V.$$

Proposition 5.8. Suppose $T \in \mathcal{L}(V, W)$. Then $\ker T \leq V$.

Proof. By 4.10, we check the conditions of a subspace:

(i) By 5.6, $T(\mathbf{0}) = \mathbf{0}$, so $\mathbf{0} \in \ker T$.

(ii) For all $v, w \in \ker T$,

$$T(v + w) = Tv + Tw = \mathbf{0} \implies v + w \in \ker T$$

so $\ker T$ is closed under addition.

(iii) For all $v \in \ker T$, $\lambda \in \mathbf{F}$,

$$T(\lambda v) = \lambda Tv = \mathbf{0} \implies \lambda v \in \ker T$$

so $\ker T$ is closed under scalar multiplication.

□

Definition 5.9 (Injectivity). Suppose $T \in \mathcal{L}(V, W)$. T is *injective* if

$$Tu = Tv \implies u = v.$$

Proposition 5.10. Suppose $T \in \mathcal{L}(V, W)$. Then T is injective if and only if $\ker T = \{\mathbf{0}\}$.

Proof.

\implies Suppose T is injective. Let $v \in \ker T$, then

$$Tv = \mathbf{0} = T(\mathbf{0}) \implies v = \mathbf{0}$$

by the injectivity of T . Hence $\ker T = \{\mathbf{0}\}$ as desired.

\impliedby Suppose $\ker T = \{\mathbf{0}\}$. Let $u, v \in V$ such that $Tu = Tv$. Then

$$T(u - v) = Tu - Tv = \mathbf{0}.$$

By definition of kernel, $u - v \in \ker T = \{\mathbf{0}\}$, so $u - v = \mathbf{0}$, which implies that $u = v$. Hence T is injective, as desired. □

Definition 5.11 (Image). Suppose $T \in \mathcal{L}(V, W)$. The *image* of T is the subset of W consisting of those

vectors that are of the form Tv for some $v \in V$:

$$\operatorname{im} T := \{Tv \mid v \in V\} \subset W.$$

Proposition 5.12. *Suppose $T \in \mathcal{L}(V, W)$. Then $\operatorname{im} T \leq W$.*

Proof.

(i) $T(\mathbf{0}) = \mathbf{0}$ implies that $\mathbf{0} \in \operatorname{im} T$.

(ii) For $w_1, w_2 \in \operatorname{im} T$, there exist $v_1, v_2 \in V$ such that $Tv_1 = w_1$ and $Tv_2 = w_2$. Then

$$w_1 + w_2 = Tv_1 + Tv_2 = T(v_1 + v_2) \in \operatorname{im} T \implies w_1 + w_2 \in \operatorname{im} T.$$

(iii) For $w \in \operatorname{im} T$ and $\lambda \in \mathbf{F}$, there exists $v \in V$ such that $Tv = w$. Then

$$\lambda w = \lambda Tv = T(\lambda v) \in \operatorname{im} T \implies \lambda w \in \operatorname{im} T.$$

□

Definition 5.13 (Surjectivity). Suppose $T \in \mathcal{L}(V, W)$. T is *surjective* if $\operatorname{im} T = W$.

Fundamental Theorem of Linear Maps

Theorem 5.14 (Fundamental theorem of linear maps). *Suppose V is finite-dimensional, $T \in \mathcal{L}(V, W)$. Then $\text{im } T$ is finite-dimensional, and*

$$\dim V = \dim \ker T + \dim \text{im } T. \quad (5.1)$$

Proof. Let $\{u_1, \dots, u_m\}$ be basis of $\ker T$, then $\dim \ker T = m$. The linearly independent list u_1, \dots, u_m can be extended to a basis

$$\{u_1, \dots, u_m, v_1, \dots, v_n\}$$

of V , thus $\dim V = m + n$. To simultaneously show that $\text{im } T$ is finite-dimensional and $\dim \text{im } T = n$, we prove that $\{Tv_1, \dots, Tv_n\}$ is a basis of $\text{im } T$. Thus we need to show that the set (i) spans $\text{im } T$, and (ii) is linearly independent.

(i) Let $v \in V$. Since $\{u_1, \dots, u_m, v_1, \dots, v_n\}$ spans V , we can write

$$v = a_1u_1 + \dots + a_mu_m + b_1v_1 + \dots + b_nv_n,$$

for some $a_i, b_i \in \mathbf{F}$. Applying T to both sides of the equation, and noting that $Tu_i = \mathbf{0}$ since $u_i \in \ker T$,

$$\begin{aligned} Tv &= T(a_1u_1 + \dots + a_mu_m + b_1v_1 + \dots + b_nv_n) \\ &= a_1 \underbrace{Tu_1}_{\mathbf{0}} + \dots + a_m \underbrace{Tu_m}_{\mathbf{0}} + b_1Tv_1 + \dots + b_nv_n \\ &= b_1Tv_1 + \dots + b_nv_n \in \text{im } T. \end{aligned}$$

Since every element of $\text{im } T$ can be expressed as a linear combination of Tv_1, \dots, Tv_n , we have that $\{Tv_1, \dots, Tv_n\}$ spans $\text{im } T$.

Moreover, since there exists a set of vectors that spans $\text{im } T$, $\text{im } T$ is finite-dimensional.

(ii) Suppose there exist $c_1, \dots, c_n \in \mathbf{F}$ such that

$$c_1Tv_1 + \dots + c_nTv_n = \mathbf{0}.$$

Then

$$T(c_1v_1 + \dots + c_nv_n) = T(\mathbf{0}) = \mathbf{0},$$

which implies $c_1v_1 + \dots + c_nv_n \in \ker T$. Since $\{u_1, \dots, u_m\}$ is a spanning set of $\ker T$, we can write

$$c_1v_1 + \dots + c_nv_n = d_1u_1 + \dots + d_mu_m$$

for some $d_i \in \mathbf{F}$, or

$$c_1v_1 + \dots + c_nv_n - d_1u_1 - \dots - d_mu_m = \mathbf{0}.$$

Since $u_1, \dots, u_m, v_1, \dots, v_n$ are linearly independent, $c_i = d_i = 0$. Since $c_i = 0$, $\{Tv_1, \dots, Tv_n\}$ is linearly independent.

□

We now show that no linear map from a finite-dimensional vector space to a “smaller” vector space can be injective, where “smaller” is measured by dimension.

Proposition 5.15. *Suppose V and W are finite-dimensional vector spaces, $\dim V > \dim W$. Then there does not exist $T \in \mathcal{L}(V, W)$ such that T is injective.*

Proof. Since W is finite-dimensional and $\text{im } T \leq W$, by 4.37, we have that $\dim \text{im } T \leq \dim W$.

Let $T \in \mathcal{L}(V, W)$. Then

$$\dim \ker T = \dim V - \dim \text{im } T \tag{1}$$

$$\geq \dim V - \dim W \tag{2}$$

$$> 0$$

where (1) follows from the fundamental theorem of linear maps, (2) follows from the above claim.

Since $\dim \ker T > 0$. This means that $\ker T$ contains some $v \in V \setminus \{0\}$. Since $\ker T \neq \{0\}$, T is not injective. \square

The next result shows that no linear map from a finite-dimensional vector space to a “bigger” vector space can be surjective, where “bigger” is also measured by dimension.

Proposition 5.16. *Suppose V and W are finite-dimensional vector spaces, $\dim V < \dim W$. Then there does not exist $T \in \mathcal{L}(V, W)$ such that T is surjective.*

Proof. Let $T \in \mathcal{L}(V, W)$. Then

$$\dim \text{im } T = \dim V - \dim \ker T \tag{1}$$

$$\leq \dim V \tag{2}$$

$$< \dim W,$$

where (1) follows from the fundamental theorem of linear maps, (2) follows since the dimension of a vector space is non-negative so $\dim \ker T \geq 0$.

Since $\dim \text{im } T < \dim W$, $\text{im } T \neq W$ so T is not surjective. \square

Example 5.17 (Homogeneous system of linear equations). Consider the homogeneous system of linear equations

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &= 0 \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n &= 0 \\ &\vdots \\ a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n &= 0 \end{aligned} \tag{*}$$

where $a_{ij} \in \mathbf{F}$.

Define $T: \mathbf{F}^n \rightarrow \mathbf{F}^m$ by

$$T(x_1, \dots, x_n) = \left(\sum_{i=1}^n a_{1i}x_i, \dots, \sum_{i=1}^n a_{mi}x_i \right).$$

The solution set of (*) is given by

$$\ker T = \left\{ (x_1, \dots, x_n) \in \mathbf{F}^n \mid \sum_{i=1}^n a_{1i}x_i = 0, \dots, \sum_{i=1}^n a_{mi}x_i = 0 \right\}.$$

Proposition. *A homogeneous system of linear equations with more variables than equations has non-zero solutions.*

Proof. If $n > m$, then

$$\begin{aligned} \dim \mathbf{F}^n > \dim \mathbf{F}^m &\implies T \text{ is not injective} \\ &\implies \ker T \neq \{\mathbf{0}\} \\ &\implies (*) \text{ has non-zero solutions} \end{aligned}$$

□

Proposition. *A system of linear equations with more equations than variables has no solution for some choice of the constant terms.*

Proof. If $n < m$, then

$$\begin{aligned} \dim \mathbf{F}^n < \dim \mathbf{F}^m &\implies T \text{ is not surjective} \\ &\implies \exists (c_1, \dots, c_m) \in \mathbf{F}^m, \forall (x_1, \dots, x_n) \in \mathbf{F}^n, T(x_1, \dots, x_n) \neq (c_1, \dots, c_m) \end{aligned}$$

Thus the choice of constant terms (c_1, \dots, c_m) is such that the system of linear equations

$$\begin{aligned} a_{11}x_1 + \dots + a_{1n}x_n &= c_1 \\ &\vdots \\ a_{m1}x_1 + \dots + a_{mn}x_n &= c_m \end{aligned}$$

has no solutions (x_1, \dots, x_n) .

□

§5.3 Matrices

Representing a Linear Map by a Matrix

Definition 5.18 (Matrix). Suppose $m, n \in \mathbb{N}$. An $m \times n$ **matrix** A is a rectangular array with m rows and n columns:

$$A = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix}$$

where $a_{ij} \in \mathbf{F}$ denotes the entry in row i , column j . We also denote $A = (a_{ij})_{m \times n}$, and drop the subscript if there is no ambiguity.

Notation. i is used for indexing across the m rows, j is used for indexing across the n columns.

Notation. $\mathcal{M}_{m \times n}(\mathbf{F})$ denotes the set of $m \times n$ matrices with entries in \mathbf{F} .

As we will soon see, matrices provide an efficient method of recording the values of Tv_j 's in terms of a basis of W .

Definition 5.19 (Matrix of linear map). Suppose $T \in \mathcal{L}(V, W)$, $\mathcal{V} = \{v_1, \dots, v_n\}$ is a basis of V , $\mathcal{W} = \{w_1, \dots, w_m\}$ is a basis of W . The matrix of T with respect to these bases is the $m \times n$ matrix $\mathcal{M}(T)$, whose entries a_{ij} are defined by

$$Tv_j = \sum_{i=1}^m a_{ij} w_i.$$

That is, the j -th column of $\mathcal{M}(T)$ consists of the scalars a_{1j}, \dots, a_{mj} needed to write Tv_j as a linear combination of the bases of W .

Notation. If the bases of V and W are not clear from the context, we adopt the notation $\mathcal{M}(T; \mathcal{V}, \mathcal{W})$.

Addition and Scalar Multiplication of Matrices

Definition 5.20 (Matrix operations).

- (i) Addition: the sum of two matrices of the same size is the matrix obtained by adding corresponding entries in the matrices:

$$\begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix} + \begin{pmatrix} c_{11} & \cdots & c_{1n} \\ \vdots & & \vdots \\ c_{m1} & \cdots & c_{mn} \end{pmatrix} = \begin{pmatrix} a_{11} + c_{11} & \cdots & a_{1n} + c_{1n} \\ \vdots & & \vdots \\ a_{m1} + c_{m1} & \cdots & a_{mn} + c_{mn} \end{pmatrix}.$$

- (ii) Scalar multiplication: the product of a scalar and a matrix is the matrix obtained by multiplying each entry in the matrix by the scalar:

$$\lambda \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix} = \begin{pmatrix} \lambda a_{11} & \cdots & \lambda a_{1n} \\ \vdots & & \vdots \\ \lambda a_{m1} & \cdots & \lambda a_{mn} \end{pmatrix}.$$

Lemma 5.21. Suppose $S, T \in \mathcal{L}(V, W)$. Then

- (i) $\mathcal{M}(S + T) = \mathcal{M}(S) + \mathcal{M}(T)$;
(ii) $\mathcal{M}(\lambda T) = \lambda \mathcal{M}(T)$ for $\lambda \in \mathbf{F}$.

Proof. Suppose $S, T \in \mathcal{L}(V, W)$, $\{v_1, \dots, v_n\}$ is a basis of V , $\{w_1, \dots, w_m\}$ is a basis of W .

- (i) By definition, $\mathcal{M}(S)$ is the matrix whose entries a_{ij} are defined by

$$Sv_j = \sum_{i=1}^m a_{ij}w_i.$$

Similarly, $\mathcal{M}(T)$ is the matrix whose entries b_{ij} are defined by

$$Tv_j = \sum_{i=1}^m b_{ij}w_i.$$

$\mathcal{M}(S + T)$ is the matrix whose entries c_{ij} are defined by

$$\begin{aligned}(S + T)v_j &= \sum_{i=1}^m c_{ij}w_i \\ Sv_j + Tv_j &= \sum_{i=1}^m c_{ij}w_i \\ \sum_{i=1}^m a_{ij}w_i + \sum_{i=1}^m b_{ij}w_i &= \sum_{i=1}^m c_{ij}w_i \\ \sum_{i=1}^m (a_{ij} + b_{ij})w_i &= \sum_{i=1}^m c_{ij}w_i \\ a_{ij} + b_{ij} &= c_{ij}.\end{aligned}$$

(ii) By definition, $\mathcal{M}(T)$ is the matrix whose entries a_{ij} are defined by

$$Tv_j = \sum_{i=1}^m a_{ij}w_i.$$

Then for $\lambda \in \mathbf{F}$, $\mathcal{M}(\lambda T)$ is the matrix whose entries b_{ij} are defined by

$$\begin{aligned}\lambda Tv_j &= \sum_{i=1}^m b_{ij}w_i \\ \lambda \sum_{i=1}^m a_{ij}w_i &= \sum_{i=1}^m b_{ij}w_i \\ \lambda a_{ij} &= b_{ij}.\end{aligned}$$

□

Lemma 5.22. *With addition and scalar multiplication defined as above, $\mathcal{M}_{m \times n}(\mathbf{F})$ is a vector space of dimension mn .*

Proof. The verification that $\mathcal{M}_{m \times n}(\mathbf{F})$ is a vector space is left to the reader. Note that the additive identity of $\mathcal{M}_{m \times n}(\mathbf{F})$ is the $m \times n$ matrix all of whose entries equal 0.

The reader should also verify that the list of distinct $m \times n$ matrices that have 0 in all entries except for a 1 in one entry is a basis of $\mathcal{M}_{m \times n}(\mathbf{F})$. There are mn such matrices, so the dimension of $\mathcal{M}_{m \times n}(\mathbf{F})$ equals mn . □

Matrix Multiplication

Note that we define the product of two matrices only when the number of columns of the first matrix equals the number of rows of the second matrix.

Definition 5.23 (Matrix multiplication). Suppose $A = (a_{ij})_{m \times n}$, $B = (b_{ij})_{n \times p}$. Then

$$AB = \left(\sum_{k=1}^n a_{ik} b_{kj} \right)_{m \times p}.$$

Remark. Thus the entry in row j , column k of AB is computed by taking row j of A and column k of B , multiplying together corresponding entries, and then summing.

In the next result, we assume that the same basis of V is used in considering $T \in \mathcal{L}(U, V)$ and $S \in \mathcal{L}(V, W)$, the same basis of W is used in considering $S \in \mathcal{L}(V, W)$ and $ST \in \mathcal{L}(U, W)$, and the same basis of U is used in considering $T \in \mathcal{L}(U, V)$ and $ST \in \mathcal{L}(U, W)$.

Lemma 5.24 (Matrix of product of linear maps). If $T \in \mathcal{L}(U, V)$ and $S \in \mathcal{L}(V, W)$, then $\mathcal{M}(ST) = \mathcal{M}(S)\mathcal{M}(T)$.

Proof. Suppose $\{v_1, \dots, v_n\}$ is a basis of V , $\{w_1, \dots, w_m\}$ is a basis of W , $\{u_1, \dots, u_p\}$ is a basis of U .

Let $\mathcal{M}(S) = A$, $\mathcal{M}(T) = B$. For $j = 1, \dots, p$, we have

$$\begin{aligned} (ST)u_j &= S(Tu_j) \\ &= S \left(\sum_{k=1}^n b_{kj} v_k \right) \\ &= \sum_{k=1}^n b_{kj} S v_k \\ &= \sum_{k=1}^n b_{kj} \left(\sum_{i=1}^m a_{ik} w_i \right) \\ &= \sum_{i=1}^m \left(\sum_{k=1}^n a_{ik} b_{kj} \right) w_i. \end{aligned}$$

□

Notation. $A_{i,\cdot}$ denotes the row vector corresponding to the i -th row of A ; $A_{\cdot,j}$ denotes the column vector corresponding to the j -th column of A .

Lemma 5.25. Suppose $A = (a_{ij})_{m \times n}$, $B = (b_{ij})_{n \times p}$. Then $AB = (c_{ij})_{m \times p}$, where

$$c_{ij} = A_{i,\cdot} B_{\cdot,j} \quad (i = 1, \dots, m, j = 1, \dots, p).$$

That is, the entry in row i , column j of AB equals (row i of A) times (column j of B).

Proof. We have

$$A_{i,\cdot}B_{\cdot,j} = \begin{pmatrix} a_{i1} & \cdots & a_{in} \end{pmatrix} \begin{pmatrix} b_{1j} \\ \vdots \\ b_{nj} \end{pmatrix} = \sum_{k=1}^n a_{ik}b_{kj} = c_{ij}.$$

□

Lemma 5.26. Suppose $A = (a_{ij})_{m \times n}$, $B = (b_{ij})_{n \times p}$. Then

$$(AB)_{\cdot,j} = AB_{\cdot,j} \quad (j = 1, \dots, p).$$

That is, column j of AB equals A times column j of B .

Proof. Using the previous result,

$$AB_{\cdot,j} = \begin{pmatrix} A_{1,\cdot}B_{\cdot,j} \\ \vdots \\ A_{n,\cdot}B_{\cdot,j} \end{pmatrix} = \begin{pmatrix} c_{1j} \\ \vdots \\ c_{nj} \end{pmatrix} = (AB)_{\cdot,j}$$

□

Lemma 5.27 (Linear combination of columns). Suppose $A = (a_{ij})_{m \times n}$, $b = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix}$. Then

$$Ab = b_1A_{\cdot,1} + \cdots + b_nA_{\cdot,n}.$$

That is, Ab is a linear combination of the columns of A , with the scalars that multiply the columns coming from b .

Proof.

$$\begin{aligned}
 Ab &= \begin{pmatrix} a_{11}b_1 + \cdots + a_{1n}b_n \\ \vdots \\ a_{m1}b_1 + \cdots + a_{mn}b_n \end{pmatrix} \\
 &= \begin{pmatrix} a_{11}b_1 \\ \vdots \\ a_{m1}b_1 \end{pmatrix} + \cdots + \begin{pmatrix} a_{1n}b_n \\ \vdots \\ a_{mn}b_n \end{pmatrix} \\
 &= b_1 \begin{pmatrix} a_{11} \\ \vdots \\ a_{m1} \end{pmatrix} + \cdots + b_n \begin{pmatrix} a_{1n} \\ \vdots \\ a_{mn} \end{pmatrix} \\
 &= b_1 A_{\cdot,1} + \cdots + b_n A_{\cdot,n}.
 \end{aligned}$$

□

The following result states that matrix multiplication can be expressed as linear combinations of columns or rows.

Proposition 5.28. *Suppose $C = (c_{ij})_{m \times c}$, $R = (r_{jk})_{c \times n}$. Then*

- (i) *Columns: for $k = 1, \dots, n$, $(CR)_{\cdot,k}$ is a linear combination of $C_{\cdot,1}, \dots, C_{\cdot,c}$, with coefficients coming from $R_{\cdot,k}$.*
- (ii) *Rows: for $i = 1, \dots, m$, $(CR)_{i,\cdot}$ is a linear combination of $R_{1,\cdot}, \dots, R_{c,\cdot}$, with coefficients coming from $C_{i,\cdot}$.*

Proof.

(i)

(ii)

□

Rank of a Matrix

Definition 5.29. Suppose $A \in \mathcal{M}_{m \times n}(\mathbf{F})$. Then the *row space* of A is the span of its rows, and the *column space* of A is the span of its columns:

$$\begin{aligned} \text{Row}(A) &:= \text{span}(A_{i,\cdot} \mid 1 \leq i \leq m), \\ \text{Col}(A) &:= \text{span}(A_{\cdot,j} \mid 1 \leq j \leq n). \end{aligned}$$

The *row rank* and *column rank* of A are defined as

$$\begin{aligned} r(A) &:= \dim \text{Row}(A), \\ c(A) &:= \dim \text{Col}(A). \end{aligned}$$

Definition 5.30 (Transpose). Suppose $A = (a_{ij})_{m \times n}$. Then the *transpose* of A is the matrix $A^T = (b_{ij})_{n \times m}$, whose entries are defined by

$$b_{ij} = a_{ji}.$$

Proposition 5.31 (Properties of transpose). Suppose $A, B \in \mathcal{M}_{m \times n}(\mathbf{F})$, $C \in \mathcal{M}_{n \times p}(\mathbf{F})$. Then

- (i) $(A + B)^T = A^T + B^T$;
- (ii) $(\lambda A)^T = \lambda A^T$ for $\lambda \in \mathbf{F}$;
- (iii) $(AC)^T = C^T A^T$.

Lemma 5.32 (Column-row factorisation). Suppose $A \in \mathcal{M}_{m \times n}(\mathbf{F})$, $c(A) \geq 1$. Then there exist $C \in \mathcal{M}_{m \times c(A)}(\mathbf{F})$, $R \in \mathcal{M}_{c(A) \times n}(\mathbf{F})$ such that $A = CR$.

Proof. We prove by construction, i.e. construct the required matrices C and R .

Each column of A is a $m \times 1$ matrix. The set of columns of A

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

can be reduced to a basis of $\text{Col}(A)$, which has length $c(A)$, by the definition of column rank. The $c(A)$ columns in this basis can be put together to form a $m \times c(A)$ matrix, which we call C .

For $k = 1, \dots, n$, the k -th column of A is a linear combination of the columns of C . Make the coefficients of this linear combination into column k of a $c \times n$ matrix, which we call R . By , it follows that $A = CR$. \square

Lemma 5.33 (Column rank equals row rank). *The column rank of a matrix equals to its row rank.*

Proof. Suppose $A \in \mathcal{M}_{m \times n}(\mathbf{F})$. Let $A = CR$ be the column-row factorisation of A , where $C \in \mathcal{M}_{m \times c(A)}(\mathbf{F})$, $R \in \mathcal{M}_{c(A) \times n}(\mathbf{F})$.

\square

Since column rank equals row rank, we can dispense with the terms “column rank” and “row rank”, and just use the simpler term “rank”.

Definition 5.34 (Rank). The *rank* of a matrix A is defined as

$$\text{rank } A := r(A) = c(A).$$

§5.4 Invertibility and Isomorphism

Invertibility

Notation. $I_V \in \mathcal{L}(V)$ denotes the identity map on V :

$$Iv = v \quad (\forall v \in V)$$

The subscript is omitted if there is no ambiguity.

Definition 5.35 (Invertibility). $T \in \mathcal{L}(V, W)$ is **invertible** if there exists $S \in \mathcal{L}(W, V)$ such that $ST = I_V, TS = I_W$; S is known as an *inverse* of T .

Lemma 5.36 (Uniqueness of inverse). *The inverse of an invertible linear map is unique.*

Proof. Suppose $T \in \mathcal{L}(V, W)$ is invertible, $S_1, S_2 \in \mathcal{L}(W, V)$ are inverses of T . Then

$$S_1 = S_1 I_W = S_1 (T S_2) = (S_1 T) S_2 = I_V S_2 = S_2.$$

Thus $S_1 = S_2$. □

Since the inverse is unique, we can give it a notation.

Notation. If T is invertible, then its inverse is denoted by T^{-1} .

The following result is useful in determining if a linear map is invertible.

Lemma 5.37 (Invertibility criterion). *Suppose $T \in \mathcal{L}(V, W)$.*

(i) T is invertible $\iff T$ is injective and surjective.

(ii) If $\dim V = \dim W$, T is invertible $\iff T$ is injective $\iff T$ is surjective.

Proof.

- (i) \implies Suppose $T \in \mathcal{L}(V, W)$ is invertible, which has inverse T^{-1} . Suppose $Tu = Tv$. Applying T^{-1} to both sides of the equation gives

$$u = T^{-1}Tu = T^{-1}Tv = v$$

so T is injective.

We now show T is surjective. Let $w \in W$. Then $w = T(T^{-1}w)$, which shows that $w \in \text{im } T$, so $\text{im } T = W$. Hence T is surjective.

\impliedby Suppose T is injective and surjective.

Define $S \in \mathcal{L}(W, V)$ such that for each $w \in W$, $S(w)$ is the unique element of V such that $T(S(w)) = w$ (we can do this due to injectivity and surjectivity). Then we have that $T(ST)v = (TS)Tv = Tv$ and thus $STv = v$ so $ST = I$. It is easy to show that S is a linear map.

(ii) It suffices to only prove T is injective $\iff T$ is surjective. Then apply the previous result.

\implies Suppose T is injective. Then $\dim \ker T = 0$. By the fundamental theorem of linear maps,

$$\begin{aligned}\dim \operatorname{im} T &= \dim V - \dim \ker T \\ &= \dim V \\ &= \dim W\end{aligned}$$

which implies that T is surjective.

\impliedby Suppose T is surjective, then $\dim \operatorname{im} T = \dim W$. By the fundamental theorem of linear maps,

$$\begin{aligned}\dim \ker T &= \dim V - \dim \operatorname{im} T \\ &= \dim V - \dim W \\ &= 0\end{aligned}$$

which implies that T is injective.

□

Corollary 5.38. *Suppose V and W are finite-dimensional, $\dim V = \dim W$, $S \in \mathcal{L}(W, V)$, $T \in \mathcal{L}(V, W)$. Then $ST = I$ if and only if $TS = I$.*

Proof.

\implies Suppose $ST = I$. Let $v \in \ker T$. Then

$$v = Iv = (ST)v = S(Tv) = S(\mathbf{0}) = \mathbf{0} \implies \ker T = \{\mathbf{0}\}$$

so T is injective. Since $\dim V = \dim W$, by the previous result, T is invertible.

Since $ST = I$, then

$$S = STT^{-1} = IT^{-1} = T^{-1}$$

so $TS = TT^{-1} = I$, as desired.

\impliedby Similar to above; reverse the roles of S and T (and V and W) to show that if $TS = I$ then $ST = I$. □

Isomorphism

Definition 5.39 (Isomorphism). An *isomorphism* is an invertible linear map. V and W are *isomorphic*, denoted by $V \cong W$, if there exists an isomorphism $T \in \mathcal{L}(V, W)$.

The following result shows that we need to look at only at the dimension to determine whether two vector spaces are isomorphic.

Lemma 5.40. *Suppose V and W are finite-dimensional. Then*

$$V \cong W \iff \dim V = \dim W.$$

Proof.

\implies Suppose $V \cong W$, then there exists an isomorphism $T \in \mathcal{L}(V, W)$, which is invertible, so T is both injective and surjective, thus $\ker T = \{0\}$ and $\text{im } T = W$, implying $\dim \ker T = 0$ and $\dim \text{im } T = \dim W$.

By the fundamental theorem of linear maps,

$$\begin{aligned} \dim V &= \dim \ker T + \dim \text{im } T \\ &= 0 + \dim W = \dim W. \end{aligned}$$

\impliedby Suppose V and W are finite-dimensional, $\dim V = \dim W = n$. Let $\{v_1, \dots, v_n\}$ be a basis of V , $\{w_1, \dots, w_n\}$ be a basis of W .

It suffices to construct an surjective $T \in \mathcal{L}(V, W)$. By the linear map lemma, there exists a linear map $T \in \mathcal{L}(V, W)$ such that

$$Tv_i = w_i \quad (i = 1, \dots, n)$$

Let $w \in W$. Then there exist $a_i \in \mathbf{F}$ such that $w = a_1 w_1 + \dots + a_n w_n$. Then

$$\begin{aligned} T(a_1 v_1 + \dots + a_n v_n) &= w \implies w \in \text{im } T \\ &\implies W = \text{im } T \\ &\implies T \text{ is surjective} \\ &\implies T \text{ is invertible.} \end{aligned}$$

□

Proposition 5.41. *Suppose $\{v_1, \dots, v_n\}$ is a basis of V , $\{w_1, \dots, w_m\}$ is a basis of W . Then*

$$\mathcal{L}(V, W) \cong \mathcal{M}_{m \times n}(\mathbf{F}).$$

Proof. We claim that \mathcal{M} is an isomorphism between $\mathcal{L}(V, W)$ and $\mathcal{M}_{m \times n}(\mathbf{F})$.

We already noted that \mathcal{M} is linear. We need to prove that \mathcal{M} is (i) injective and (ii) surjective.

(i) Given $T \in \mathcal{L}(V, W)$, if $\mathcal{M}(T) = 0$, then

$$Tv_j = 0 \quad (j = 1, \dots, n)$$

Since v_1, \dots, v_n is a basis of V , this implies $T = \mathbf{0}$, so $\ker \mathcal{M} = \{\mathbf{0}\}$. Thus \mathcal{M} is injective.

(ii) Suppose $A \in \mathcal{M}_{m \times n}(\mathbf{F})$. By the linear map lemma, there exists $T \in \mathcal{L}(V, W)$ such that

$$Tv_j = \sum_{i=1}^m a_{ij}w_i \quad (j = 1, \dots, n)$$

Since $\mathcal{M}(T) = A$, $\text{im } \mathcal{M} = \mathcal{M}_{m \times n}(\mathbf{F})$ so \mathcal{M} is surjective.

□

The following is a useful corollary.

Lemma 5.42. *Suppose V and W are finite-dimensional. Then $\mathcal{L}(V, W)$ is finite-dimensional and*

$$\dim \mathcal{L}(V, W) = (\dim V)(\dim W).$$

Proof. Since $\mathcal{L}(V, W) \cong \mathcal{M}_{m \times n}(\mathbf{F})$,

$$\dim \mathcal{L}(V, W) = \dim \mathcal{M}_{m \times n}(\mathbf{F}) = mn = (\dim V)(\dim W).$$

□

Linear Maps Thought of as Matrix Multiplication

Previously we defined the matrix of a linear map. Now we define the matrix of a vector.

Definition 5.43 (Matrix of a vector). Suppose $v \in V$, $\{v_1, \dots, v_n\}$ is a basis of V . The matrix of v with respect to this basis is

$$\mathcal{M}(v) = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix}$$

where $b_1, \dots, b_n \in \mathbf{F}$ are such that

$$v = b_1 v_1 + \dots + b_n v_n.$$

Example 5.44. If $x = (x_1, \dots, x_n) \in \mathbf{F}^n$, then the matrix of the vector x with respect to the standard basis is

$$\mathcal{M}(x) = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}.$$

Lemma 5.45. Suppose $T \in \mathcal{L}(V, W)$. Let $\{v_1, \dots, v_n\}$ be a basis of V , $\{w_1, \dots, w_m\}$ be a basis of W . Then

$$\mathcal{M}(T)_{.j} = \mathcal{M}(Tv_j) \quad (j = 1, \dots, n)$$

Proof. By definition, the entries of $\mathcal{M}(T)$ are defined such that

$$Tv_j = \sum_{i=1}^m a_{ij} w_i \quad (j = 1, \dots, n)$$

Then since $Tv_j \in W$, by definition, the matrix of Tv_j with respect to the basis $\{w_1, \dots, w_m\}$ is

$$\mathcal{M}(Tv_j) = \begin{pmatrix} a_{1j} \\ \vdots \\ a_{mj} \end{pmatrix}$$

which is precisely the j -th column of $\mathcal{M}(T)_{.j}$, for $j = 1, \dots, n$. □

The following result shows that linear maps act like matrix multiplication.

Lemma 5.46. Suppose $T \in \mathcal{L}(V, W)$. Let $\{v_1, \dots, v_n\}$ be a basis of V , $\{w_1, \dots, w_m\}$ be a basis of W . Let $v \in V$, then

$$\mathcal{M}(Tv) = \mathcal{M}(T)\mathcal{M}(v).$$

Proof. Suppose $v = b_1v_1 + \cdots + b_nv_n$ for some $b_1, \dots, b_n \in \mathbf{F}$. Then

$$\begin{aligned} \mathcal{M}(Tv) &= \mathcal{M}(T(b_1v_1 + \cdots + b_nv_n)) \\ &= b_1\mathcal{M}(Tv_1) + \cdots + b_n\mathcal{M}(Tv_n) \\ &= b_1\mathcal{M}(T)_{\cdot,1} + \cdots + b_n\mathcal{M}(T)_{\cdot,n} \\ &= \begin{pmatrix} \mathcal{M}(T)_{\cdot,1} & \cdots & \mathcal{M}(T)_{\cdot,n} \end{pmatrix} \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} \\ &= \mathcal{M}(T)\mathcal{M}(v). \end{aligned}$$

□

Notice that no bases are in sight in the statement of the next result. Although $\mathcal{M}(T)$ in the next result depends on a choice of bases of V and W , the next result shows that the column rank of $\mathcal{M}(T)$ is the same for all such choices (because $\text{im } T$ does not depend on a choice of basis).

Proposition 5.47. *Suppose V and W are finite-dimensional, $T \in \mathcal{L}(V, W)$. Then*

$$\dim \ker T = \text{rank } \mathcal{M}(T).$$

Proof. Suppose $\{v_1, \dots, v_n\}$ is a basis of V , $\{w_1, \dots, w_m\}$ is a basis of W .

The linear map that takes $w \in W$ to $\mathcal{M}(w)$ is an isomorphism from W to $\mathcal{M}_{m \times 1}(\mathbf{F})$ (consisting of $m \times 1$ column vectors).

The restriction of this isomorphism to $\text{im } T$ [which equals $\text{span}(Tv_1, \dots, Tv_n)$] is an isomorphism from $\text{im } T$ to $\text{span}(\mathcal{M}(Tv_1), \dots, \mathcal{M}(Tv_n))$. For $j = 1, \dots, n$, the $m \times 1$ matrix $\mathcal{M}(Tv_j)$ equals column j of $\mathcal{M}(T)$. Thus

$$\dim \ker T = \text{rank } \mathcal{M}(T),$$

as desired. □

Change of Basis

Definition 5.48 (Identity matrix). For $n \in \mathbb{N}$, the $n \times n$ *identity matrix* is

$$I_n = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix}.$$

Remark. Note that the symbol I is used to denote both the identity operator and the identity matrix. The context indicates which meaning of I is intended. For example, consider the equation $\mathcal{M}(I) = I$; on LHS I denotes the identity operator, and on RHS I denotes the identity matrix.

The next result justifies the name “identity matrix”.

Lemma 5.49. *Suppose $A \in \mathcal{M}_{n \times n}(\mathbf{F})$. Then $AI_n = I_n A = A$.*

Proof. Exercise. □

Definition 5.50 (Invertible matrix). $A \in \mathcal{M}_{n \times n}(\mathbf{F})$ is called *invertible* if there exists $B \in \mathcal{M}_{n \times n}(\mathbf{F})$ such that $AB = BA = I$; we call B an *inverse* of A

Lemma 5.51 (Uniqueness of inverse). *Suppose A is an invertible square matrix. Then there exists a unique matrix B such that $AB = BA = I$.*

Proof. Let B and C be inverses of A . Then

$$B = BI = BAC = IC = C.$$

□

Since the inverse of a matrix is unique, we can give it a notation.

Notation. The inverse of a matrix A is denoted by A^{-1} .

Lemma 5.52.

- (i) *Suppose A is an invertible square matrix. Then $(A^{-1})^{-1} = A$.*
- (ii) *Suppose A and C are invertible square matrices of the same size. Then AC is invertible, and $(AC)^{-1} = C^{-1}A^{-1}$.*

Proof.

(i) We have

$$A^{-1}A = AA^{-1} = I,$$

so the inverse of A^{-1} is A .

(ii) We have

$$\begin{aligned} (AC)(C^{-1}A^{-1}) &= A(CC^{-1})A^{-1} \\ &= AIA^{-1} \\ &= AA^{-1} \\ &= I, \end{aligned}$$

and similarly $(C^{-1}A^{-1})(AC) = I$.

□

Proposition 5.53 (Matrix of product of linear maps). *Suppose $T \in \mathcal{L}(U, V)$, $S \in \mathcal{L}(V, W)$. Let $\mathcal{U} = \{u_1, \dots, u_m\}$ be a basis of U , $\mathcal{V} = \{v_1, \dots, v_n\}$ be a basis of V , $\mathcal{W} = \{w_1, \dots, w_p\}$ be a basis of W . Then*

$$\mathcal{M}(ST; \mathcal{U}, \mathcal{W}) = \mathcal{M}(S; \mathcal{V}, \mathcal{W}) \mathcal{M}(T; \mathcal{U}, \mathcal{V}).$$

Proof. Refer to previous section. Now we are just being more explicit about the bases involved. □

Corollary 5.54. *Suppose that $\mathcal{U} = \{u_1, \dots, u_n\}$ and $\mathcal{V} = \{v_1, \dots, v_n\}$ are bases of V . Then the matrices*

$$\mathcal{M}(I; \mathcal{U}, \mathcal{V}) \quad \text{and} \quad \mathcal{M}(I; \mathcal{V}, \mathcal{U})$$

are invertible, and each is the inverse of the other.

Theorem 5.55 (Change-of-basis formula). *Suppose $T \in \mathcal{L}(V)$. Let $\mathcal{U} = \{u_1, \dots, u_n\}$ and $\mathcal{V} = \{v_1, \dots, v_n\}$ be bases of V . Let*

$$A = \mathcal{M}(T; \mathcal{U}), \quad B = \mathcal{M}(T; \mathcal{V}),$$

and $C = \mathcal{M}(I; \mathcal{U}, \mathcal{V})$. Then

$$A = C^{-1}BC. \tag{5.2}$$

Proof. Note that

$$\begin{aligned} \mathcal{M}(T; \mathcal{U}, \mathcal{V}) &= \underbrace{\mathcal{M}(T; \mathcal{V})}_B \underbrace{\mathcal{M}(I; \mathcal{U}, \mathcal{V})}_C \\ &= \underbrace{\mathcal{M}(I; \mathcal{U}, \mathcal{V})}_C \underbrace{\mathcal{M}(T; \mathcal{U})}_A \end{aligned}$$

Hence $BC = CA$, and the desired result follows. □

The next result states that the matrix of inverse equals inverse of matrix.

Lemma 5.56. *Suppose $\{v_1, \dots, v_n\}$ is a basis of V , $T \in \mathcal{L}(V)$ is invertible. Then*

$$\mathcal{M}(T^{-1}) = (\mathcal{M}(T))^{-1},$$

where both matrices are with respect to the basis $\{v_1, \dots, v_n\}$.

Proof. We have that

$$\mathcal{M}(T^{-1})\mathcal{M}(T) = \mathcal{M}(T^{-1}T) = \mathcal{M}(I) = I.$$

□

§5.5 Products and Quotients of Vector Spaces

Products of Vector Spaces

Definition 5.57 (Product). Suppose V_1, \dots, V_n are vector spaces over \mathbf{F} . The **product** $V_1 \times \dots \times V_n$ is defined by

$$V_1 \times \dots \times V_n := \{(v_1, \dots, v_n) \mid v_i \in V_i\}.$$

Remark. This is analogous to the Cartesian product of sets.

Lemma 5.58. $V_1 \times \dots \times V_n$ is a vector space over \mathbf{F} , with addition and scalar multiplication defined by

$$\begin{aligned} (u_1, \dots, u_n) + (v_1, \dots, v_n) &= (u_1 + v_1, \dots, u_n + v_n) \\ \lambda(v_1, \dots, v_n) &= (\lambda v_1, \dots, \lambda v_n) \end{aligned}$$

The following result shows that the dimension of a product is the sum of dimensions.

Proposition 5.59. Suppose V_1, \dots, V_n are finite-dimensional. Then $V_1 \times \dots \times V_n$ is finite-dimensional, and

$$\dim(V_1 \times \dots \times V_n) = \dim V_1 + \dots + \dim V_n.$$

Proof. For each V_k ($k = 1, \dots, n$), choose a basis:

$$\mathcal{B}_k = \{e_{k1}, \dots, e_{k \dim V_k}\}.$$

For each basis vector of each V_k , consider the set consisting of elements of $V_1 \times \dots \times V_n$ that equal the basis vector in the k -th slot and 0 in the other slots:

$$\mathcal{B} = \{(0, \dots, \underbrace{e_{ki}}_{k\text{-th slot}}, \dots, 0) \mid 1 \leq i \leq \dim V_k, 1 \leq k \leq n\}.$$

We want to show that \mathcal{B} is a basis of $V_1 \times \dots \times V_n$. Thus we need to show that it is (i) a spanning set, and (ii) linearly independent.

(i) Let $(v_1, \dots, v_n) \in V_1 \times \dots \times V_n$. For $k = 1, \dots, n$, since \mathcal{B}_k is a basis for V_k , we can write

$$v_k = \sum_{i=1}^{\dim V_k} a_{ki} e_{ki}.$$

for some $a_{k1}, \dots, a_{k \dim V_k} \in \mathbf{F}$. Then

$$\begin{aligned} (v_1, \dots, v_n) &= \sum_{k=1}^n (0, \dots, v_k, \dots, 0) \\ &= \sum_{k=1}^n \left(0, \dots, \sum_{i=1}^{\dim V_k} a_{ki} e_{ki}, \dots, 0 \right) \\ &= \sum_{k=1}^n \sum_{i=1}^{\dim V_k} a_{ki} (0, \dots, e_{ki}, \dots, 0) \end{aligned}$$

which is a linear combination of vectors in \mathcal{B} . Hence \mathcal{B} spans $V_1 \times \dots \times V_n$.

(ii) Suppose there exist $a_{ki} \in \mathbf{F}$ such that

$$\begin{aligned} \sum_{k=1}^n \sum_{i=1}^{\dim V_k} a_{ki} (0, \dots, e_{ki}, \dots, 0) &= \mathbf{0} \\ \sum_{k=1}^n \left(0, \dots, \sum_{i=1}^{\dim V_k} a_{ki} e_{ki}, \dots, 0 \right) &= \mathbf{0} \\ \left(\sum_{i=1}^{\dim V_1} a_{1i} e_{1i}, \sum_{i=1}^{\dim V_2} a_{2i} e_{2i}, \dots, \sum_{i=1}^{\dim V_n} a_{ni} e_{ni} \right) &= \mathbf{0} \end{aligned}$$

so for $k = 1, \dots, n$,

$$\sum_{i=1}^{\dim V_k} a_{ki} e_{ki} = \mathbf{0}.$$

By the linear independence of vectors in \mathcal{B}_k , we have that

$$a_{k1} = \dots = a_{k \dim V_k} = 0$$

for $k = 1, \dots, n$.

Hence

$$\begin{aligned} \dim(V_1 \times \dots \times V_n) &= |\mathcal{B}| \\ &= |\mathcal{B}_1| + \dots + |\mathcal{B}_n| \\ &= \dim V_1 + \dots + \dim V_n. \end{aligned}$$

□

Products are also related to direct sums, by the following result.

Proposition 5.60. *Suppose that $V_1, \dots, V_n \leq V$. Define a linear map*

$$\begin{aligned} \Gamma : V_1 \times \dots \times V_n &\rightarrow V_1 + \dots + V_n \\ (v_1, \dots, v_n) &\mapsto v_1 + \dots + v_n \end{aligned}$$

Then $V_1 + \cdots + V_n$ is a direct sum if and only if Γ is injective.

Proof.

(i) \iff (ii) Suppose $V_1 + \cdots + V_n$ is a direct sum. Let $(v_1, \dots, v_n) \in \ker \Gamma$. Then

$$\Gamma(v_1, \dots, v_n) = \mathbf{0}$$

$$v_1 + \cdots + v_n = \mathbf{0}$$

$$v_1 = \cdots = v_n = \mathbf{0}$$

so $(v_1, \dots, v_n) = \mathbf{0}$. Hence $\ker \Gamma = \mathbf{0}$, thus Γ is injective.

(ii) \iff (i) Similar to the above proof. □

The next result says that a sum is a direct sum if and only if dimensions add up.

Proposition 5.61. *Suppose V is finite-dimensional, $V_1, \dots, V_n \leq V$. Then $V_1 + \cdots + V_n$ is a direct sum if and only if*

$$\dim(V_1 + \cdots + V_n) = \dim V_1 + \cdots + \dim V_n.$$

Proof. The map Γ defined in the previous result is surjective. Thus by the fundamental theorem of linear maps, Γ is injective if and only if

$$\dim(V_1 + \cdots + V_n) = \dim(V_1 \times \cdots \times V_n).$$

Then use the previous two results above. □

Quotient Spaces

Definition 5.62 (Coset). Suppose $v \in V, U \subset V$. Then $v + U$ is called a *coset* of U , defined by

$$v + U := \{v + u \mid u \in U\}.$$

Definition 5.63 (Quotient space). Suppose $U \leq V$. Then the *quotient space* V/U is the set of cosets of U :

$$V/U := \{v + U \mid v \in V\}.$$

Example 5.64. If $U = \{(x, 2x) \in \mathbb{R}^2 \mid x \in \mathbb{R}\}$, then \mathbb{R}^2/U is the set of lines in \mathbb{R}^2 that have gradient of 2.

It is obvious that two cosets of a subspace are equal or disjoint. We shall now prove this.

Lemma 5.65. Suppose $U \leq V$, and $v, w \in V$. Then

$$v - w \in U \iff v + U = w + U \iff (v + U) \cap (w + U) = \emptyset.$$

Proof. First suppose $v - w \in U$. If $u \in U$, then

$$v + u = w + ((v - w) + u) \in w + U.$$

Thus $v + U \subset w + U$. Similarly, $w + U \subset v + U$. Thus $v + U = w + U$, completing the proof that $v - w \in U$ implies $v + U = w + U$.

The equation $v + U = w + U$ implies that $(v + U) \cap (w + U) \neq \emptyset$.

Now suppose $(v + U) \cap (w + U) \neq \emptyset$. Thus there exist $u_1, u_2 \in U$ such that

$$v + u_1 = w + u_2.$$

Thus $v - w = u_2 - u_1$. Hence $v - w \in U$, showing that $(v + U) \cap (w + U) \neq \emptyset$ implies $v - w \in U$, which completes the proof. \square

We can define a vector space structure on V/U .

Lemma 5.66. Suppose $U \leq V$. Then V/U is a vector space, with addition and scalar multiplication defined by

$$\begin{aligned} (v + U) + (w + U) &= (v + w) + U \\ \lambda(v + U) &= (\lambda v) + U \end{aligned}$$

for all $v, w \in V, \lambda \in \mathbf{F}$.

Proof. We first need to show that addition and scalar multiplication are well-defined.

Addition Suppose $v_1, v_2, w_1, w_2 \in V$ are such that

$$v_1 + U = v_2 + U, \quad w_1 + U = w_2 + U.$$

By 5.65,

$$v_1 - v_2 \in U, \quad w_1 - w_2 \in U.$$

Since $U \leq V$, U is closed under addition, so $(v_1 - v_2) + (w_1 - w_2) \in U$. Thus $(v_1 + w_1) - (v_2 + w_2) \in U$. Using 5.65 again, we see that

$$(v_1 + w_1) + U = (v_2 + w_2) + U,$$

as desired. Hence addition on V/U is well-defined.

Scalar multiplication Suppose $v_1, v_2 \in V$ are such that $v_1 + U = v_2 + U$, suppose $\lambda \in \mathbf{F}$.

Since $U \leq V$, U is closed under scalar multiplication, so $\lambda(v_1 - v_2) \in U$. Thus $\lambda v_1 - \lambda v_2 \in U$. By 5.65,

$$(\lambda v_1) + U = (\lambda v_2) + U.$$

Hence scalar multiplication on V/U is well-defined.

The verification that addition and scalar multiplication make V/U into a vector space is straightforward and is left to the reader. Note that the additive identity of V/U is $0 + U$ (which equals U) and that the additive inverse of $v + U$ is $(-v) + U$. \square

Definition 5.67 (Quotient map). Suppose $U \leq V$. The *quotient map* is the map

$$\begin{aligned} \pi : V &\rightarrow V/U \\ v &\mapsto v + U \end{aligned}$$

for all $v \in V$.

We check that the quotient map is a linear map: let $v, w \in V$, $\lambda \in \mathbf{F}$,

$$(i) \quad \pi(v) + \pi(w) = (v + U) + (w + U) = (v + w) + U = \pi(v + w).$$

$$(ii) \quad \pi(\lambda v) = (\lambda v) + U = \lambda(v + U) = \lambda(\pi v).$$

Proposition 5.68 (Dimension of quotient space). Suppose V is finite-dimensional, $U \leq V$. Then

$$\dim V/U = \dim V - \dim U.$$

Idea. Since dimensions are involved, think of the fundamental theorem of linear maps.

Proof. Let the quotient map $\pi : V \rightarrow V/U$.

- Let $v \in V$. Then

$$\begin{aligned} v \in \ker \pi &\iff \pi(v) = \mathbf{0} + U = U \\ &\iff v + U = U \quad [\text{by 5.65}] \\ &\iff v \in U \end{aligned}$$

so $\ker \pi = U$.

- The definition of π implies $\text{im } \pi = V/U$.

By the fundamental theorem of linear maps,

$$\begin{aligned}\dim V &= \dim \ker \pi + \dim \operatorname{im} \pi \\ &= \dim U + \dim V/U\end{aligned}$$

which gives the desired result. \square

Each linear map T on V induces a linear map \tilde{T} on $V/\ker T$, as defined below.

Definition 5.69. Suppose $T \in \mathcal{L}(V, W)$. Define

$$\begin{aligned}\tilde{T} : V/\ker T &\rightarrow W \\ v + \ker T &\mapsto Tv\end{aligned}$$

We first show that \tilde{T} is well-defined.

Proof. Suppose $u, v \in V$ are such that

$$u + \ker T = v + \ker T.$$

By 5.65, $u - v \in \ker T$. Thus $T(u - v) = \mathbf{0}$, so $Tu = Tv$. \square

We then check that \tilde{T} is a linear map from $V/\ker T$ to W .

The next result shows that we can think of \tilde{T} as a modified version of T , with a domain that produces an injective map.

Proposition 5.70. Suppose $T \in \mathcal{L}(V, W)$. Then

- (i) $\tilde{T} \circ \pi = T$, where π is the quotient map of V onto $V/\ker T$;
- (ii) \tilde{T} is injective;
- (iii) $\operatorname{im} \tilde{T} = \operatorname{im} T$.

Proof.

(i) Let $v \in V$. Then

$$(\tilde{T} \circ \pi)(v) = \tilde{T}(\pi(v)) = \tilde{T}(v + \ker T) = Tv.$$

(ii) Let $v + \ker T \in \ker \tilde{T}$. Then

$$\begin{aligned}\tilde{T}(v + \ker T) = \mathbf{0} &\implies Tv = \mathbf{0} \\ &\implies v \in \ker T \\ &\implies v + \ker T = \ker T\end{aligned}$$

so $\ker \tilde{T} = \{\mathbf{0} + \ker T\}$.

(iii) The definition of \tilde{T} shows that $\operatorname{im} \tilde{T} = \operatorname{im} T$.

□

Theorem 5.71 (First isomorphism theorem). *Suppose $T \in \mathcal{L}(V, W)$ is an isomorphism. Then*

$$V/\ker T \cong \operatorname{im} T. \quad (5.3)$$

Proof. Now (ii) and (iii) imply that if we think of \tilde{T} as mapping into $\operatorname{im} T$, then \tilde{T} is an isomorphism from $V/\ker T$ onto $\operatorname{im} T$. □

Theorem 5.72 (Second isomorphism theorem).

Theorem 5.73 (Third isomorphism theorem). *$U \subset V \subset W$, then*

$$W/V \cong (W/U)/(V/U). \quad (5.4)$$

§5.6 Duality

Dual Space and Dual Map

Linear maps into the scalar field \mathbf{F} play a special role in linear algebra, so they get a special name.

Definition 5.74. A *linear functional* on V is a linear map from V to \mathbf{F} . (That is, a linear functional is an element of $\mathcal{L}(V, \mathbf{F})$.)

The *dual space* of V is the vector space of linear functionals on V ; that is, $V^* := \mathcal{L}(V, \mathbf{F})$.

Lemma 5.75 (Dimension of dual space). *Suppose V is finite-dimensional. Then V^* is finite-dimensional, and*

$$\dim V^* = \dim V.$$

Proof. By 5.42,

$$\dim V^* := \dim \mathcal{L}(V, \mathbf{F}) = (\dim V)(\dim \mathbf{F}) = \dim V.$$

□

Definition 5.76 (Dual basis). Let $\{v_1, \dots, v_n\}$ be a basis of V . Then the *dual basis* of $\{v_1, \dots, v_n\}$ is

$$\{\phi_1, \dots, \phi_n\} \subset V^*,$$

where each ϕ_i is the linear functional on V such that

$$\phi_i(v_j) = \begin{cases} 1 & (i = j) \\ 0 & (i \neq j) \end{cases}$$

Example 5.77 (Dual basis of the standard basis of \mathbf{F}^n). Fix a positive integer n . For $i = 1, \dots, n$, define ϕ_i to be the linear functional on \mathbf{F}^n that selects the i -th coordinate of a vector in \mathbf{F}^n :

$$\phi_i(x_1, \dots, x_n) = x_i$$

for each $(x_1, \dots, x_n) \in \mathbf{F}^n$.

Let e_1, \dots, e_n be the standard basis of \mathbf{F}^n . Then

$$\phi_i(e_j) = \begin{cases} 1 & (i = j) \\ 0 & (i \neq j) \end{cases}$$

Thus ϕ_1, \dots, ϕ_n is the dual basis of the standard basis e_1, \dots, e_n of \mathbf{F}^n .

The next result shows that the dual basis of a basis of V consists of the linear functionals on V that give the coefficients for expressing a vector in V as a linear combination of the basis vectors.

Proposition 5.78. *Suppose $\{v_1, \dots, v_n\}$ is a basis of V , and $\{\phi_1, \dots, \phi_n\}$ is the dual basis. Then for each $v \in V$,*

$$v = \phi_1(v)v_1 + \dots + \phi_n(v)v_n.$$

Proof. Let $v \in V$. Since v_1, \dots, v_n is a basis of V , there exist $c_1, \dots, c_n \in \mathbf{F}$ such that

$$v = c_1v_1 + \dots + c_nv_n.$$

For $i = 1, \dots, n$, applying ϕ_i to both sides of the equation above gives

$$\phi_i(v) = c_i.$$

□

The next result shows that the dual basis is a basis of the dual space. Thus the terminology “dual basis” is justified.

Lemma 5.79. *Suppose V is finite-dimensional. Then the dual basis of a basis of V is a basis of V^* .*

Proof. Suppose v_1, \dots, v_n is a basis of V . Let ϕ_1, \dots, ϕ_n denote the dual basis.

Claim. ϕ_1, \dots, ϕ_n is linearly independent in V^* .

Suppose $a_1, \dots, a_n \in \mathbf{F}$ are such that

$$a_1\phi_1 + \dots + a_n\phi_n = 0. \tag{1}$$

Now for each $i = 1, \dots, n$,

$$(a_1\phi_1 + \dots + a_n\phi_n)(v_i) = a_i.$$

Thus (1) shows that $a_1 = \dots = a_n = 0$. Hence ϕ_1, \dots, ϕ_n is linearly independent.

Since ϕ_1, \dots, ϕ_n is a linearly independent set in V of length $\dim V$, we can conclude that ϕ_1, \dots, ϕ_n is a basis of V . □

Definition 5.80 (Dual map). Suppose $T \in \mathcal{L}(V, W)$. The **dual map** of T is the linear map

$$\begin{aligned} T^* : W^* &\rightarrow V^* \\ \phi &\mapsto \phi \circ T \end{aligned}$$

We check that $T^* \in \mathcal{L}(W^*, V^*)$: let $\phi, \psi \in W^*$, $\lambda \in \mathbf{F}$,

$$(i) \quad T^*(\phi + \psi) = (\phi + \psi) \circ T = \phi \circ T + \psi \circ T = T^*(\phi) + T^*(\psi)$$

$$(ii) \quad T^*(\lambda\phi) = (\lambda\phi) \circ T = \lambda(\phi \circ T) = \lambda(T^*(\phi))$$

Lemma 5.81 (Algebraic properties of dual map). *Suppose $T \in \mathcal{L}(V, W)$. Then*

$$(i) \quad (S + T)^* = S^* + T^* \text{ for all } S \in \mathcal{L}(V, W)$$

$$(ii) (\lambda T)^* = \lambda T^* \text{ for all } \lambda \in \mathbf{F}$$

$$(iii) (ST)^* = T^*S^* \text{ for all } S \in \mathcal{L}(V, W)$$

Proof.

(i)

(ii)

(iii) Let $\phi \in U^*$. Then

$$(ST)^*(\phi) = \phi \circ (ST) = (\phi \circ S) \circ T = T^*(\phi \circ S) = T^*(S^*(\phi)) = (T^*S^*)(\phi).$$

□

(i) and (ii) imply that the function that takes T to T^* is a linear map from $\mathcal{L}(V, W)$ to $\mathcal{L}(W^*, V^*)$.

Kernel and Image of Dual of Linear Map

The goal of this section is to describe $\ker T^*$ and $\operatorname{im} T^*$ in terms of $\operatorname{im} T$ and $\ker T$. To do this, we will need the next definition.

Definition 5.82 (Annihilator). For $U \subset V$, the *annihilator* of U is defined by

$$U^0 := \{\phi \in V^* \mid \phi(u) = \mathbf{0}, \forall u \in U\}.$$

We check that $U^0 \leq V$:

- (i) Note that $0 \in U^0$ (here 0 is the zero linear functional on V) because the zero linear functional applied to every vector in U equals $\mathbf{0} \in \mathbf{F}$.
- (ii) Suppose $\phi, \psi \in U^0$. Thus $\phi, \psi \in V^*$ and $\phi(u) = \psi(u) = \mathbf{0}$ for every $u \in U$.

Let $u \in U$, then

$$(\phi + \psi)(u) = \phi(u) + \psi(u) = \mathbf{0} + \mathbf{0} = \mathbf{0}.$$

Thus $\phi + \psi \in U^0$, so U^0 is closed under addition.

- (iii) Suppose $\phi \in U^0$, $\lambda \in \mathbf{F}$, let $u \in U$, then

$$\phi(\lambda u) = \lambda \phi(u) = \mathbf{0}$$

so $\lambda \phi \in U^0$, so U^0 is closed under scalar multiplication.

Proposition 5.83 (Dimension of annihilator). Suppose V is finite-dimensional, and $U \leq V$. Then

$$\dim U^0 = \dim V - \dim U.$$

Proof.

□

The following are conditions for the annihilator to equal $\{\mathbf{0}\}$ or the whole space.

Proposition 5.84. Suppose V is finite-dimensional, and $U \leq V$. Then

- (i) $U^0 = \{\mathbf{0}\} \iff U = V$;
- (ii) $U^0 = V^* \iff U = \{\mathbf{0}\}$.

The following result concerns $\ker T^*$.

Proposition 5.85. Suppose V and W are finite-dimensional, $T \in \mathcal{L}(V, W)$. Then

- (i) $\ker T^* = (\operatorname{im} T)^0$;
- (ii) $\dim \ker T^* = \dim \ker T + \dim W - \dim V$.

The next result can be useful because sometimes it is easier to verify that T^* is injective than to show directly that T is surjective.

Proposition 5.86. *Suppose V and W are finite-dimensional, $T \in \mathcal{L}(V, W)$. Then*

$$T \text{ is surjective} \iff T^* \text{ is injective.}$$

The following result concerns $\text{im } T^*$.

Proposition 5.87. *Suppose V and W finite-dimensional, $T \in \mathcal{L}(V, W)$. Then*

(i) $\dim \text{im } T^* = \dim \text{im } T$;

(ii) $\dim T^* = (\ker T)^0$.

Proposition 5.88. *Suppose V and W are finite-dimensional, $T \in \mathcal{L}(V, W)$. Then*

$$T \text{ is injective} \iff T^* \text{ is surjective.}$$

Matrix of Dual of Linear Map

Proposition 5.89. *Suppose V and W are finite-dimensional, $T \in \mathcal{L}(V, W)$. Then*

$$\mathcal{M}(T^*) = (\mathcal{M}(T))^t.$$

Exercises

Exercise 5.1 ([Axl24] 3A). Suppose $b, c \in \mathbb{R}$. Define $T: \mathbb{R}^3 \rightarrow \mathbb{R}^2$ by

$$T(x, y, z) = (2x - 4y + 3z + b, 6x + cxyz).$$

Show that T is linear if and only if $b = c = 0$.

Exercise 5.2 ([Axl24] 3A Q11). Suppose V is finite-dimensional, $T \in \mathcal{L}(V)$. Prove that T is a scalar multiple of the identity if and only if $ST = TS$ for all $S \in \mathcal{L}(V)$.

Exercise 5.3 ([Axl24] 3B Q9). Suppose $T \in \mathcal{L}(V, W)$ is injective, $\{v_1, \dots, v_n\}$ is linearly independent in V . Prove that $\{Tv_1, \dots, Tv_n\}$ is linearly independent in W .

Solution. Suppose there exist $a_i \in \mathbf{F}$ such that

$$\begin{aligned} a_1Tv_1 + \dots + a_nTv_n &= \mathbf{0} \\ \implies T(a_1v_1 + \dots + a_nv_n) &= \mathbf{0} \\ \implies a_1v_1 + \dots + a_nv_n &\in \ker T \end{aligned}$$

Since T is injective,

$$\ker T = \{\mathbf{0}\} \implies a_1v_1 + \dots + a_nv_n = \mathbf{0} \implies a_1 = \dots = a_n = 0$$

since $\{v_1, \dots, v_n\}$ is linearly independent. □

Exercise 5.4 ([Axl24] 3B Q11). Suppose that V is finite-dimensional, $T \in \mathcal{L}(V, W)$. Prove that there exists $U \leq V$ such that

$$U \cap \ker T = \{\mathbf{0}\} \quad \text{and} \quad \text{im } T = T(U).$$

Solution. □

Exercise 5.5 ([Axl24] 3B Q19). Suppose W is finite-dimensional, $T \in \mathcal{L}(V, W)$. Prove that T is injective if and only if there exists $S \in \mathcal{L}(W, V)$ such that ST is the identity operator on V .

Solution. □

Exercise 5.6 ([Axl24] 3B Q20). Suppose W is finite-dimensional, $T \in \mathcal{L}(V, W)$. Prove that T is surjective if and only if there exists $S \in \mathcal{L}(W, V)$ such that TS is the identity operator on W .

Exercise 5.7 ([Axl24] 3B 22). Suppose U, V are finite-dimensional, $S \in \mathcal{L}(V, W)$, $T \in \mathcal{L}(U, V)$. Prove that

$$\dim \ker ST \leq \dim \ker S + \dim \ker T.$$

Solution. □

Exercise 5.8 ([Axl24] 3D). Suppose $T \in \mathcal{L}(V, W)$ is invertible. Show that T^{-1} is invertible and

$$\left(T^{-1}\right)^{-1} = T.$$

Solution. T^{-1} is invertible because there exists T such that $TT^{-1} = T^{-1}T = I$. So

$$T^{-1}T = TT^{-1} = I$$

thus $(T^{-1})^{-1} = T$. □

3C Q15,16,17

3D Q11,12,17,22,23,24

Exercise 5.9 ([Ax124] 3D). Suppose $T \in \mathcal{L}(U, V)$ and $S \in \mathcal{L}(V, W)$ are both invertible linear maps. Prove that $ST \in \mathcal{L}(U, W)$ is invertible and that $(ST)^{-1} = T^{-1}S^{-1}$.

Solution.

$$(ST)(T^{-1}S^{-1}) = S(TT^{-1})S^{-1} = I = T^{-1}S^{-1}ST.$$

□

Exercise 5.10 ([Ax124] 3D). Suppose V is finite-dimensional and $T \in \mathcal{L}(V, W)$. Prove that the following are equivalent:

- (i) T is invertible;
- (ii) $\{Tv_1, \dots, Tv_n\}$ is a basis of V for every basis $\{v_1, \dots, v_n\}$ of V ;
- (iii) $\{Tv_1, \dots, Tv_n\}$ is a basis of V for some basis $\{v_1, \dots, v_n\}$ of V .

Solution.

(i) \implies (ii) It only suffices to prove linear independence. We can show this

$$a_1Tv_1 + \dots + a_nTv_n = 0 \iff a_1v_1 + \dots + a_nv_n = 0$$

since T is injective and thus the only solution is all a_i are identically zero.

(ii) \implies (iii) Trivial.

(iii) \implies (i) By the linear map lemma, there exists $S \in \mathcal{L}(V)$ such that $S(Tv_i) = v_i$ for all i . Such S is the inverse of T (one can verify) and thus T is invertible. □

Exercise 5.11 ([Ax124] 3E Q3). Suppose V_1, \dots, V_m are vector spaces. Prove that

$$\mathcal{L}(V_1 \times \dots \times V_m, W) \cong \mathcal{L}(V_1, W) \times \dots \times \mathcal{L}(V_m, W).$$

Exercise 5.12 ([Ax124] 3E Q4). Suppose V_1, \dots, V_m are vector spaces. Prove that

$$\mathcal{L}(V, W_1 \times \dots \times W_m) \cong \mathcal{L}(V, W_1) \times \dots \times \mathcal{L}(V, W_m).$$

Exercise 5.13 ([Ax124] 3E Q5). For a positive integer m , define V^m by

$$V^m = \underbrace{V \times \dots \times V}_{m \text{ times}}.$$

Prove that $V^m \cong \mathcal{L}(\mathbf{F}^m, V)$.

Exercise 5.14 ([Axl24] 3E Q6). Suppose that $v, x \in V$ and $U, W \leq V$ are such that $v + U = x + W$. Prove that $U = W$.

Exercise 5.15 ([Axl24] 3E Q12, Barycentric coordinates). Suppose $v_1, \dots, v_m \in V$. Let

$$A = \{\lambda_1 v_1 + \dots + \lambda_m v_m \mid \lambda_i \in \mathbf{F}, \lambda_1 + \dots + \lambda_m = 1\}.$$

- (i) Prove that A is a coset of some subspace of V .
- (ii) Prove that if B is a coset of some subspace of V , and $\{v_1, \dots, v_m\} \subset B$, then $A \subset B$.
- (iii) Prove that A is a coset of some subspace of V , where $\dim V < m$.

Exercise 5.16 ([Axl24] 3E Q13). Suppose $U \leq V$, and V/U is finite-dimensional. Prove that $V \cong U \times (V/U)$.

Solution.

$$\dim V = \dim U + (\dim V - \dim U) = \dim U + \dim(V/U).$$

□

Exercise 5.17 ([Axl24] 3E Q14). Suppose $U, W \leq V$ such that $V = U \oplus W$. Suppose w_1, \dots, w_m is a basis of W . Prove that $w_1 + U, \dots, w_m + U$ is a basis of V/U .

Exercise 5.18 ([Axl24] 3E Q15).

Exercise 5.19 ([Axl24] 3E Q16). Suppose $\phi \in \mathcal{L}(V, \mathbf{F})$ and $\phi \neq 0$. Prove that $\dim V / \ker \phi = 1$.

Exercise 5.20 ([Axl24] 3E Q18).

Exercise 5.21 ([Axl24] 3E Q19). Suppose $T \in \mathcal{L}(V, W)$ and $U \leq V$. Let π denote the quotient map from V to V/U . Prove that there exists $S \in \mathcal{L}(V/U, W)$ such that

$$T = S \circ \pi \iff U \subset \ker T.$$

6 Polynomials

§6.1 Definitions

Definition 6.1 (Polynomial). $p : \mathbf{F} \rightarrow \mathbf{F}$ is a *polynomial* with coefficients in \mathbf{F} if there exist $a_i \in \mathbf{F}$ such that

$$p(z) = a_0 + a_1z + \cdots + a_nz^n \quad (z \in \mathbf{F})$$

Notation. The set of polynomials with coefficients in \mathbf{F} is denoted by $\mathbf{F}[z]$.

Lemma 6.2. *With the usual operations of addition and scalar multiplication, $\mathbf{F}[z]$ is a vector space over \mathbf{F} .*

Hence $\mathbf{F}[z]$ is a subspace of $\mathbf{F}^{\mathbf{F}}$ (vector space of functions from \mathbf{F} to \mathbf{F}).

Definition 6.3 (Degree). A polynomial $p \in \mathbf{F}[z]$ has *degree* n , denoted by $\deg p = n$, if there exist scalars $a_0, a_1, \dots, a_n \in \mathbf{F}$ with $a_n \neq 0$ such that $p(z) = a_0 + a_1z + \cdots + a_nz^n$ for all $z \in \mathbf{F}$.

Notation. For non-negative integer n , $\mathbf{F}_n[z]$ denotes the set of polynomials with coefficients in \mathbf{F} and degree at most n .

Lemma 6.4. *For non-negative integer n , $\mathbf{F}_n[z]$ is finite-dimensional.*

Proof. $\mathbf{F}_n[z] = \text{span}(1, z, z^2, \dots, z^n)$ [here we slightly abuse notation by letting z^k denote a function]. \square

Lemma 6.5. *$\mathbf{F}[z]$ is infinite-dimensional.*

Proof. Consider any list of elements of $\mathbf{F}[z]$. Let n denote the highest degree of the polynomials in this list. Then every polynomial in the span of this list has degree at most n . Thus z^{n+1} is not in the span of our list. Hence no list spans $\mathbf{F}[z]$. Thus $\mathbf{F}[z]$ is infinite-dimensional. \square

§6.2 Zeros of Polynomials

Definition 6.6 (Zero of polynomial). $\lambda \in \mathbf{F}$ is called a **zero** of a polynomial $p \in \mathbf{F}[z]$ if

$$p(\lambda) = 0.$$

Lemma 6.7 (Factor theorem). Suppose $n \in \mathbb{N}$, $p \in \mathbf{F}_n[z]$. Suppose $\lambda \in \mathbf{F}$, then $p(\lambda) = 0$ if and only if there exists $q \in \mathbf{F}_{n-1}[z]$ such that

$$p(z) = (z - \lambda)q(z) \quad (\forall z \in \mathbf{F})$$

Proof.

\Rightarrow Suppose $p(\lambda) = 0$. Let $a_0, a_1, \dots, a_n \in \mathbf{F}$ be such that

$$p(z) = a_n z^n + \dots + a_1 z + a_0 \quad (\forall z \in \mathbf{F})$$

Then for all $z \in \mathbf{F}$,

$$\begin{aligned} p(z) &= p(z) - p(\lambda) \\ &= (a_n z^n + \dots + a_1 z + a_0) - (a_n \lambda^n + \dots + a_1 \lambda + a_0) \\ &= a_n (z^n - \lambda^n) + \dots + a_1 (z - \lambda). \end{aligned}$$

Note that for each $k = 1, \dots, n$, we can factorise

$$z^k - \lambda^k = (z - \lambda) \left(z^{k-1} + z^{k-2} \lambda + \dots + \lambda^{k-1} \right).$$

Thus p equals $z - \lambda$ times some polynomial of degree $n - 1$, as desired.

\Leftarrow Now suppose that there exists a polynomial $q \in \mathbf{F}[z]$ such that

$$p(z) = (z - \lambda)q(z) \quad (\forall z \in \mathbf{F})$$

Then

$$p(\lambda) = (\lambda - \lambda)q(\lambda) = 0,$$

as desired. □

Now we can prove that the degree of a polynomials determines how many zeros it has.

Proposition 6.8. Suppose $n \in \mathbb{N}$, $p \in \mathbf{F}_n[z]$. Then p has at most n zeros in \mathbf{F} .

Proof. Prove by induction on n .

The desired result holds for $n = 1$ because if $a_1 \neq 0$ then the polynomial $a_0 + a_1 z$ has only one zero (which equals $-\frac{a_0}{a_1}$).

Now assume the desired result holds for $n - 1$. If p has no zeros in \mathbf{F} , then the desired result holds and we are

done. Thus suppose p has a zero $\lambda \in \mathbf{F}$. By 6.7, there exists $q \in \mathbf{F}[z]$ of degree $n - 1$ such that

$$p(z) = (z - \lambda)q(z) \quad (\forall z \in \mathbf{F})$$

By the induction hypothesis, q has at most $n - 1$ zeros in \mathbf{F} . The equation above shows that the zeros of p in \mathbf{F} are exactly the zeros of q in \mathbf{F} along with λ . Thus p has at most n zeros in \mathbf{F} . \square

The result above implies that the coefficients of a polynomial are uniquely determined (because if a polynomial had two different sets of coefficients, then subtracting the two representations of the polynomial would give a polynomial with some nonzero coefficients but infinitely many zeros). In particular, the degree of a polynomial is uniquely defined.

§6.3 Division Algorithm for Polynomials

Proposition 6.9 (Division algorithm). *Suppose $p, s \in \mathbf{F}[z]$, $s \neq 0$. Then there exists unique polynomials $q, r \in \mathbf{F}[z]$, where $\deg r < \deg s$, such that*

$$p = sq + r.$$

Proof. Let $n = \deg p$, $m = \deg s$. If $n < m$, take $q = 0$ and $r = p$ to get the desired equation.

Now assume that $n \geq m$. The set

$$S = \{1, z, \dots, z^{m-1}, s, zs, \dots, z^{n-m}s\}$$

is linearly independent in $\mathbf{F}[z]$ because each polynomial in S has a different degree. Also, S has length $n + 1$, which equals $\dim \mathbf{F}[z]$. Hence S is a basis of $\mathbf{F}[z]$.

Since $p \in \mathbf{F}[z]$ and S is a basis of $\mathbf{F}[z]$, there exist unique constants $a_0, a_1, \dots, a_{m-1} \in \mathbf{F}$ and $b_0, b_1, \dots, b_{n-m} \in \mathbf{F}$ such that

$$\begin{aligned} p &= a_0 + a_1z + \dots + a_{m-1}z^{m-1} + b_0s + b_1zs + \dots + b_{n-m}z^{n-m}s \\ &= \underbrace{a_0 + a_1z + \dots + a_{m-1}z^{m-1}}_r + s \left(\underbrace{b_0 + b_1z + \dots + b_{n-m}z^{n-m}}_q \right). \end{aligned}$$

With r and q as defined above, we see that p can be written as $p = sq + r$ with $\deg r < \deg s$, as desired.

The uniqueness of $q, r \in \mathbf{F}[z]$ satisfying these conditions follows from the uniqueness of the constants $a_0, a_1, \dots, a_{m-1} \in \mathbf{F}$ and $b_0, b_1, \dots, b_{n-m} \in \mathbf{F}$. \square

§6.4 Factorisation of Polynomials over \mathbb{C}

Theorem 6.10 (Fundamental theorem of algebra, first version). *Every non-constant polynomial with complex coefficients has a zero in \mathbb{C} .*

Theorem 6.11 (Fundamental theorem of algebra). *If $p \in \mathbb{C}[z]$ is a non-constant polynomial, then p has a unique factorisation (except for the order of the factors) of the form*

$$p(z) = c(z - \lambda_1) \cdots (z - \lambda_n),$$

where $c, \lambda_1, \dots, \lambda_n \in \mathbb{C}$.

§6.5 Factorisation of Polynomials over \mathbb{R}

A polynomial with real coefficients may have no real zeros. For example, the polynomial $x^2 + 1$ has no real zeros.

To obtain a factorisation theorem over \mathbb{R} , we will use our factorisation theorem over \mathbb{C} . We begin with the next result.

Proposition 6.12. *Suppose $p \in \mathbb{C}[z]$ is a polynomial with real coefficients. If $\lambda \in \mathbb{C}$ is a zero of p , then so is the conjugate $\bar{\lambda}$.*

Proof. Let

$$p(z) = a_0 + a_1z + \cdots + a_nz^n,$$

where $a_0, \dots, a_n \in \mathbb{R}$. Suppose $\lambda \in \mathbb{C}$ is a zero of p , then

$$a_0 + a_1\lambda + \cdots + a_n\lambda^n = 0.$$

Taking the complex conjugate on both sides of the equation gives

$$a_0 + a_1\bar{\lambda} + \cdots + a_n\bar{\lambda}^n = 0.$$

Hence $\bar{\lambda}$ is a zero of p . □

We want a factorisation theorem for polynomials with real coefficients. We begin with the following result.

Remark. Think about the quadratic formula in connection with the result below.

Lemma 6.13 (Factorisation of quadratic polynomial). *Suppose $b, c \in \mathbb{R}$. Then there is a polynomial factorisation of the form*

$$x^2 + bx + c = (x - \lambda_1)(x - \lambda_2)$$

with $\lambda_1, \lambda_2 \in \mathbb{R}$ if and only if $b^2 \geq 4c$.

Proof. Completing the square gives

$$x^2 + bx + c = \left(x + \frac{b}{2}\right)^2 + \left(c - \frac{b^2}{4}\right). \quad (1)$$

\Rightarrow We prove the contrapositive. Suppose $b^2 < 4c$, then the RHS of (1) is positive for every $x \in \mathbb{R}$. Hence the polynomial $x^2 + bx + c$ has no real zeros and thus cannot be factored in the form $(x - \lambda_1)(x - \lambda_2)$ with $\lambda_1, \lambda_2 \in \mathbb{R}$.

\Leftarrow Suppose $b^2 \geq 4c$. Then there is a real number d such that $d^2 = \frac{b^2}{4} - c$. Then (1) can be written as

$$\begin{aligned} x^2 + bx + c &= \left(x + \frac{b}{2}\right)^2 - d^2 \\ &= \left(x + \frac{b}{2} + d\right) \left(x + \frac{b}{2} - d\right), \end{aligned}$$

which gives the desired factorisation. □

Theorem 6.14 (Factorisation of polynomial over \mathbb{R}). *Suppose $p \in \mathbb{R}[x]$ is a non-constant polynomial. Then p has a unique factorisation (except for the order of the factors) of the form*

$$p(x) = c(x - \lambda_1) \cdots (x - \lambda_n)(x^2 + b_1x + c_1) \cdots (x^2 + b_Nx + c_N),$$

where $c, \lambda_1, \dots, \lambda_n, b_1, \dots, b_N, c_1, \dots, c_N \in \mathbb{R}$, with $b_k^2 < 4c_k$ for each k .

7 Eigenvalues and Eigenvectors

§7.1 Invariant Subspaces

Eigenvalues

Definition 7.1 (Operator). An *operator* is a linear map from a vector space to itself.

Definition 7.2 (Invariant subspace). Suppose $T \in \mathcal{L}(V)$. $U \leq V$ is *invariant* under T if $Tu \in U$ for all $u \in U$.

Example 7.3. Suppose $T \in \mathcal{L}(V)$. Then the following subspaces of V are all invariant under T .

- (i) The subspace $\{0\}$ is invariant under T : if $u \in \{0\}$, then $u = \mathbf{0}$ so $Tu = \mathbf{0} \in \{0\}$.
- (ii) The subspace V is invariant under T : if $u \in V$, then $Tu \in V$.
- (iii) The subspace $\ker T$ is invariant under T : if $u \in \ker T$, then $Tu = \mathbf{0}$, and hence $Tu \in \ker T$, since a subspace must contain $\mathbf{0}$.
- (iv) The subspace $\operatorname{im} T$ is invariant under T : if $u \in \operatorname{im} T$, then $Tu \in \operatorname{im} T$ by definition.

Definition 7.4 (Eigenvalue and eigenvector). Suppose $T \in \mathcal{L}(V)$. $\lambda \in \mathbf{F}$ is an *eigenvalue* of T if there exists $v \in V \setminus \{0\}$ such that $Tv = \lambda v$; we say v is an *eigenvector* of T corresponding to λ .

Lemma 7.5 (Equivalent conditions to be an eigenvalue). Suppose V is finite-dimensional, $T \in \mathcal{L}(V)$, $\lambda \in \mathbf{F}$. Then the following are equivalent:

- (i) λ is an eigenvalue of T .
- (ii) $T - \lambda I$ is not injective.
- (iii) $T - \lambda I$ is not surjective.
- (iv) $T - \lambda I$ is not invertible.

Proof.

(i) \iff (ii) $Tv = \lambda v$ is equivalent to the equation $(T - \lambda I)v = \mathbf{0}$, so $T - \lambda I$ is not injective.

(ii) \iff (iii) \iff (iv) This directly follows from 5.37. □

Proposition 7.6 (Linearly independent eigenvectors). *Suppose $T \in \mathcal{L}(V)$. Then every set of eigenvectors of T corresponding to distinct eigenvalues of T is linearly independent.*

Proof. Suppose, for a contradiction, that the desired result is false. Then there exists a smallest positive integer m such that v_1, \dots, v_m are linearly dependent eigenvectors of T corresponding to distinct eigenvalues $\lambda_1, \dots, \lambda_m$ of T . The linear dependence implies there exists $a_1, \dots, a_m \in \mathbf{F}$, none of which are 0 (because of the minimality of m) such that

$$a_1 v_1 + \cdots + a_m v_m = \mathbf{0}.$$

Applying $T - \lambda_m I$ to both sides of the equation,

$$\begin{aligned} a_1(T - \lambda_m I)v_1 + \cdots + a_{m-1}(T - \lambda_m I)v_{m-1} + a_m(T - \lambda_m I)v_m &= \mathbf{0} \\ a_1(Tv_1 - \lambda_m v_1) + \cdots + a_{m-1}(Tv_{m-1} - \lambda_m v_{m-1}) + a_m(Tv_m - \lambda_m v_m) &= \mathbf{0} \\ a_1(\lambda_1 - \lambda_m)v_1 + \cdots + a_{m-1}(\lambda_{m-1} - \lambda_m)v_{m-1} &= \mathbf{0} \end{aligned}$$

Since the eigenvalues $\lambda_1, \dots, \lambda_m$ are distinct, none of the coefficients $a_i(\lambda_i - \lambda_m)$ equal 0. Thus v_1, \dots, v_{m-1} are $m - 1$ linearly dependent eigenvectors of T corresponding to distinct eigenvalues, contradicting the minimality of m . \square

Corollary 7.7. *Suppose V is finite-dimensional. Then each operator on V has at most $\dim V$ distinct eigenvalues.*

Proof. Let $T \in \mathcal{L}(V)$. Suppose $\lambda_1, \dots, \lambda_m$ are distinct eigenvalues of T with corresponding eigenvectors v_1, \dots, v_m .

By 7.6, the eigenvectors v_1, \dots, v_m are linearly independent. Since the length of a linearly independent set is less than or equal to the length of a spanning set, we have that $m \leq \dim V$, as desired. \square

Polynomials Applied to Operators

Notation. Suppose $T \in \mathcal{L}(V)$, $n \in \mathbb{Z}^+$. $T^n \in \mathcal{L}(V)$ is defined by $T^n = \underbrace{T \cdots T}_{m \text{ times}}$. T^0 is defined to be the identity operator I on V . If T is invertible with inverse T^{-1} , then $T^{-n} \in \mathcal{L}(V)$ is defined by $T^{-n} = (T^{-1})^n$.

Having defined powers of an operator, we can now define what it means to apply a polynomial to an operator.

Definition 7.8. Suppose $T \in \mathcal{L}(V)$, $p \in \mathbf{F}[z]$ is a polynomial given by

$$p(z) = a_n z^n + \cdots + a_1 z + a_0 \quad (z \in \mathbf{F})$$

Then $p(T)$ is the operator on V defined by

$$p(T) := a_n T^n + \cdots + a_1 T + a_0.$$

If we fix an operator $T \in \mathcal{L}(V)$, then the function $\mathbf{F}[z] \rightarrow \mathcal{L}(V)$ given by $p \mapsto p(T)$ is linear:

Definition 7.9 (Product of polynomials). Suppose $p, q \in \mathbf{F}[z]$. Then $pq \in \mathbf{F}[z]$ is the polynomial defined by

$$(pq)(z) = p(z)q(z) \quad (z \in \mathbf{F})$$

Lemma 7.10. Suppose $p, q \in \mathbf{F}[z]$, $T \in \mathcal{L}(V)$. Then

$$(i) \quad (pq)(T) = p(T)q(T); \quad \text{(multiplicativity)}$$

$$(ii) \quad p(T)q(T) = q(T)p(T). \quad \text{(commutativity)}$$

This means when a product of polynomials is expanded using the distributive property, it does not matter whether the symbol is z or T .

Proof.

(i) Suppose

$$p(z) = \sum_{i=0}^m a_i z^i, \quad q(z) = \sum_{j=0}^n b_j z^j \quad (z \in \mathbf{F})$$

Then

$$\begin{aligned} (pq)(z) &= p(z)q(z) \\ &= \left(\sum_{i=0}^m a_i z^i \right) \left(\sum_{j=0}^n b_j z^j \right) \\ &= \sum_{i=0}^m \sum_{j=0}^n a_i b_j z^{i+j}. \end{aligned}$$

Thus

$$\begin{aligned} (pq)(T) &= \sum_{i=0}^m \sum_{j=0}^n a_i b_j T^{i+j} \\ &= \left(\sum_{i=0}^m a_i T^i \right) \left(\sum_{j=0}^n b_j T^j \right) \\ &= p(T)q(T). \end{aligned}$$

(ii) Using (i) twice, we have

$$p(T)q(T) = (pq)(T) = (qp)(T) = q(T)p(T)$$

since the multiplication of polynomials is commutative.

□

Proposition 7.11. *Suppose $T \in \mathcal{L}(V)$, $p \in \mathbf{F}[z]$. Then*

(i) $\ker p(T)$ is invariant under T ;

(ii) $\operatorname{im} p(T)$ is invariant under T .

Proof.

(i) Let $u \in \ker p(T)$. Then $p(T)u = \mathbf{0}$. Thus

$$(p(T))(Tu) = (p(T)T)(u) = (Tp(T))(u) = T(p(T)u) = T(\mathbf{0}) = \mathbf{0}.$$

Hence $Tu \in \ker p(T)$, so $\ker p(T)$ is invariant under T .

(ii) Let $u \in \operatorname{im} p(T)$. Then there exists $v \in V$ such that $u = p(T)v$. Thus

$$Tu = T(p(T)v) = p(T)(Tv).$$

Hence $Tu \in \operatorname{im} p(T)$, so $\operatorname{im} p(T)$ is invariant under T .

□

§7.2 The Minimal Polynomial

Existence of Eigenvalues on Complex Vector Spaces

The following is one of the most important results in linear algebra.

Theorem 7.12 (Existence of eigenvalues). *Every operator on a finite-dimensional, non-zero, complex vector space has an eigenvalue.*

Proof. Suppose V is a finite-dimensional complex vector space, $\dim V = n > 0$, $T \in \mathcal{L}(V)$. Let $v \in V \setminus \{\mathbf{0}\}$. Consider the set

$$S = \{v, Tv, T^2v, \dots, T^nv\}.$$

Since $\dim V = n$ and S has length $n + 1$, S is not linearly independent; thus there exist $a_0, \dots, a_n \in \mathbb{C}$, not all 0, such that

$$a_0v + a_1Tv + a_2T^2v + \dots + a_nT^nv = \mathbf{0},$$

which we can write as

$$p(T)v = \mathbf{0},$$

where $p(z) = a_0 + a_1z + \dots + a_nz^n$, where we pick p such that $\deg p$ is minimal.

By the fundamental theorem of algebra, there exists a root of p in \mathbb{C} ; let $\lambda \in \mathbb{C}$ be a root of p . By the factor theorem,

$$p(z) = (z - \lambda)q(z) \quad (z \in \mathbb{C}).$$

Thus

$$\begin{aligned} p(T) &= (T - \lambda I)q(T) \\ \mathbf{0} &= p(T)v = (T - \lambda I)q(T)v \\ Tq(T)v &= \lambda q(T)v \end{aligned}$$

Since p is the minimal polynomial and $\deg q < \deg p$, we must have that $q(T)v \neq \mathbf{0}$. Therefore λ is an eigenvalue of T , with corresponding eigenvector $q(T)v$. \square

Example 7.13. Note that the hypothesis in 7.12 that $\mathbf{F} = \mathbb{C}$ cannot be replaced with the hypothesis that $\mathbf{F} = \mathbb{R}$.

For instance, consider $T \in \mathcal{L}(\mathbb{R}^2)$ defined by

$$Tv = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} v. \quad (*)$$

Then

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

Notice that T is a rotation, so there is no vector that is fixed in its original direction. Hence T does not have an eigenvalue.

In contrast, consider $T \in \mathcal{L}(\mathbb{C}^2)$ defined by (*). Then

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

so i is an eigenvalue with corresponding eigenvector $\begin{pmatrix} i \\ 1 \end{pmatrix}$.

Eigenvalues and the Minimal Polynomial

A *monic polynomial* is a polynomial whose highest-degree coefficient equals 1.

The following result shows the existence, uniqueness and degree of the *minimal polynomial*.

Lemma 7.14. *Suppose V is finite-dimensional, $T \in \mathcal{L}(V)$. Then there exists a unique monic polynomial $p \in \mathbf{F}[z]$ of smallest degree such that $p(T) = 0$. Furthermore, $\deg p \leq \dim V$.*

Proof.

Existence Let $\dim V = n$. We use strong induction on n .

If $n = 0$, then T is the zero operator on V ; thus take p to be the constant polynomial 1.

Now assume that $n > 0$ and that the desired result holds for all operators on all vector spaces of smaller dimension. We want to construct a monic polynomial of smallest degree such that when applied to T gives the 0 operator.

Let $u \in V \setminus \{\mathbf{0}\}$, consider the set

$$\{u, Tu, T^2u, \dots, T^nu\}.$$

This set has length $n + 1$, so it is linearly dependent. By the linear dependence lemma, there exists a smallest positive integer $m \leq n$ such that $T^m u$ is a linear combination of $u, Tu, \dots, T^{m-1}u$; thus there exist $c_i \in \mathbf{F}$ such that

$$c_0u + c_1Tu + \dots + c_{m-1}T^{m-1}u + T^m u = \mathbf{0}.$$

Define a monic polynomial $q \in \mathbf{F}[z]$ by $q(z) = c_0 + c_1z + \dots + c_{m-1}z^{m-1} + z^m$. Then $q(T)u = \mathbf{0}$. Thus for non-negative integer k ,

$$q(T)(T^k u) = T^k(q(T)u) = T^k(\mathbf{0}) = \mathbf{0}.$$

By the linear dependence lemma, $\{u, Tu, \dots, T^{m-1}u\}$ is linearly independent. Thus the above equation implies that $\dim \ker q(T) \geq m$. Hence by the fundamental theorem of linear maps,

$$\begin{aligned} \dim \operatorname{im} q(T) &= \dim V - \dim \ker q(T) \\ &\leq \dim V - m. \end{aligned}$$

Since $\operatorname{im} q(T)$ is invariant under T , we can apply the induction hypothesis to the restriction $T|_{\operatorname{im} q(T)}$. Thus there exists a monic polynomial $s \in \mathbf{F}[z]$ with $\deg s \leq \dim V - m$ such that

$$s(T|_{\operatorname{im} q(T)}) = 0.$$

Hence for all $v \in V$ we have

$$((sq)(T))v = s(T)(q(T)v) = \mathbf{0}$$

because $q(T)v \in \operatorname{im} q(T)$ and $s(T)|_{\operatorname{im} q(T)} = s(T|_{\operatorname{im} q(T)}) = 0$. Thus sq is a monic polynomial such that $\deg sq \leq \dim V$ and $(sq)(T) = 0$, as desired.

Uniqueness Let $p \in \mathbf{F}[z]$ be a monic polynomial of smallest degree such that $p(T) = 0$; let $r \in \mathbf{F}[z]$ be a monic polynomial of same degree and $r(T) = 0$. Then $(p - r)(T) = 0$ and also $\deg(p - r) < \deg p$.

We claim that $p - r = 0$. Suppose otherwise, for a contradiction, that $p - r \neq 0$. Then divide $p - r$ by the coefficient of the highest-order term in $p - r$ to get a monic polynomial $s \in \mathbf{F}[z]$, which satisfies $s(T) = 0$ and

also $\deg s = \deg(p - r) < \deg p$, a contradiction. \square

Definition 7.15 (Minimal polynomial). Suppose V is finite-dimensional, $T \in \mathcal{L}(V)$. The *minimal polynomial* of T is the unique monic polynomial $p \in \mathbf{F}[z]$ of smallest degree such that $p(T) = 0$.

Theorem 7.16. Suppose V is finite-dimensional, $T \in \mathcal{L}(V)$.

- (i) The zeros of the minimal polynomial of T are eigenvalues of T .
- (ii) If V is a complex vector space, then the minimal polynomial of T has the form

$$(z - \lambda_1) \cdots (z - \lambda_m),$$

where λ_i are eigenvalues of T .

Proof. Let p be the minimal polynomial of T .

- (i) First suppose $\lambda \in \mathbf{F}$ is a zero of p . Then p can be written in the form

$$p(z) = (z - \lambda)q(z)$$

where q is a monic polynomial with coefficients in \mathbf{F} . Since $p(T) = 0$, we have

$$\mathbf{0} = (T - \lambda I)(q(T)v) \quad (v \in V).$$

Since $\deg p < \deg p$ and p is the minimal polynomial of T , there exists at least one $v \in V$ such that $q(T)v \neq \mathbf{0}$. The equation above thus implies that λ is an eigenvalue of T , as desired.

To prove that every eigenvalue of T is a zero of p , now suppose $\lambda \in \mathbf{F}$ is an eigenvalue of T . Thus there exists $v \in V \setminus \{\mathbf{0}\}$ such that $Tv = \lambda v$. Repeated applications of T to both sides of this equation show that $T^k v = \lambda^k v$ for every nonnegative integer k . Thus

$$p(T)v = p(\lambda)v.$$

Since p is the minimal polynomial of T , we have $p(T)v = \mathbf{0}$. Hence the equation above implies that $p(\lambda)v = \mathbf{0}$. Thus λ is a zero of p , as desired.

- (ii) To get the desired result, use (i) and the second version of the fundamental theorem of algebra. \square

The next result completely characterises the polynomials that when applied to an operator give the 0 operator.

Proposition 7.17. Suppose V is finite-dimensional, $T \in \mathcal{L}(V)$, $q \in \mathbf{F}[z]$. Then $q(T) = 0$ if and only if q is a polynomial multiple of the minimal polynomial of T .

Proof. Let p denote the minimal polynomial of T .

\implies Suppose $q(T) = 0$. By the division algorithm, there exist polynomials $s, r \in \mathbf{F}[z]$ such that

$$q + ps + r \tag{1}$$

and $\deg r < \deg p$. We have

$$0 = q(T) = p(T)s(T) + r(T) = r(T).$$

The equation above implies that $r = 0$ (otherwise, dividing r by its highest-degree coefficient would produce a monic polynomial that when applied to T gives 0; this polynomial would have a smaller degree than the minimal polynomial, which would be a contradiction). Thus (1) becomes $q = ps$, so q is a polynomial multiple of p .

$\boxed{\Leftarrow}$ Suppose q is a polynomial multiple of p . Thus $q = ps$ for some polynomial $s \in \mathbf{F}[z]$, so

$$q(T) = p(T)s(T) = 0s(T) = 0$$

as desired. □

The next result is a nice consequence of the above result.

Corollary 7.18. *Suppose V is finite-dimensional, $T \in \mathcal{L}(V)$, $U \leq V$ is invariant under T . Then the minimal polynomial of T is a polynomial multiple of the minimal polynomial of $T|_U$.*

Proof. Let p be the minimal polynomial of T . Then $p(T)v = \mathbf{0}$ for all $v \in V$. In particular,

$$p(T)u = \mathbf{0} \quad (\forall u \in U).$$

Thus $p(T|_U) = 0$. By 7.17 (applied to $T|_U$ in place of T), p is a polynomial multiple of the minimal polynomial of $T|_U$. □

The next result shows that the constant term of the minimal polynomial of an operator determines whether the operator is invertible.

Corollary 7.19. *Suppose V is finite-dimensional, $T \in \mathcal{L}(V)$. Then T is not invertible if and only if the constant term of the minimal polynomial of T is 0.*

Proof. Suppose $T \in \mathcal{L}(V)$, let p be the minimal polynomial of T . Then

$$T \text{ is not invertible} \iff 0 \text{ is an eigenvalue of } T \quad [\text{by 7.5}]$$

$$\iff 0 \text{ is a zero of } p \quad [\text{by 7.16}]$$

$$\iff \text{constant term of } p \text{ is } 0.$$

□

Eigenvalues on Odd-Dimensional Real Vector Spaces

The next result will be the key tool that we use to show that every operator on an odd-dimensional real vector space has an eigenvalue.

Lemma 7.20. *Suppose V is a finite-dimensional, real vector space. Suppose $T \in \mathcal{L}(V)$ and $b, c \in \mathbb{R}$ with $b^2 < 4c$. Then $\dim \ker(T^2 + bT + cI)$ is even.*

Proof. By 7.11, $\ker(T^2 + bT + cI)$ is invariant under T . By replacing V with $\ker(T^2 + bT + cI)$ and replacing T with T restricted to $\ker(T^2 + bT + cI)$, we can assume that $T^2 + bT + cI = 0$; we now need to prove that $\dim V$ is even.

Suppose $\lambda \in \mathbb{R}$ and $v \in V$ are such that $Tv = \lambda v$. Then

$$\mathbf{0} = (T^2 + bT + cI)v = (\lambda^2 + b\lambda + c)v = \underbrace{\left(\left(\lambda + \frac{b}{2} \right)^2 + c - \frac{b^2}{4} \right)}_{>0} v$$

so $v = \mathbf{0}$. Hence we have shown that T has no eigenvectors.

Let $U \leq V$ be invariant under T , and has the largest dimension among all subspaces of V that are invariant under T and have even dimension.

Claim. $U = V$.

If $U = V$, then we are done; otherwise assume there exists $w \in V$ such that $w \notin U$.

Let $W = \text{span}(w, Tw)$. Then W is invariant under T because $T(Tw) = -bTw - cw$. Furthermore, $\dim W = 2$ because otherwise w would be an eigenvector of T . Now

$$\dim(U + W) = \dim U + \dim W - \dim(U \cap W) = \dim U + 2,$$

where $U \cap W = \{\mathbf{0}\}$ because otherwise $U \cap W$ would be a one-dimensional subspace of V that is invariant under T (impossible because T has no eigenvectors).

Because $U + W$ is invariant under T , the equation above shows that there exists a subspace of V invariant under T of even dimension larger than $\dim U$. Thus the assumption that $U \neq V$ was incorrect. Hence V has even dimension. \square

The next result states that on odd-dimensional vector spaces, every operator has an eigenvalue. We already know this result for finite-dimensional complex vector spaces (without the odd hypothesis). Thus in the proof below, we will assume that $\mathbf{F} = \mathbb{R}$.

Proposition 7.21. *Every operator on an odd-dimensional vector space has an eigenvalue.*

Proof. Suppose V is a finite-dimensional real vector space, $\dim V = n$ is odd. Let $T \in \mathcal{L}(V)$. We will induct on n in steps of size two to show that T has an eigenvalue.

To get started, note that the desired result holds if $\dim V = 1$ because then every nonzero vector in V is an eigenvector of T .

Now suppose that $n \geq 3$ and the desired result holds for all operators on all odd-dimensional vector spaces of dimension less than n . Let p be the minimal polynomial of T .

If p is a polynomial multiple of $x - \lambda$ for some $\lambda \in \mathbb{R}$, then λ is an eigenvalue of T [by 5.27(a)] and we are done. Thus we can assume that there exist $b, c \in \mathbb{R}$ such that $b^2 < 4c$ and p is a polynomial multiple of $x^2 + bx + c$ (see 4.16).

There exists a monic polynomial $q \in \mathbb{R}[x]$ such that $p(x) = q(x)(x^2 + bx + c)$ for all $x \in \mathbb{R}$. Now

$$0 = p(T) = (q(T))(T^2 + bT + cI),$$

which means that $q(T)$ equals 0 on $\text{im}(T^2 + bT + cI)$. Because $\deg q < \deg p$ and p is the minimal polynomial of T , this implies that $\text{im}(T^2 + bT + cI) \neq V$.

By the fundamental theorem of linear maps,

$$\dim V = \dim \ker(T^2 + bT + cI) + \dim \text{im}(T^2 + bT + cI).$$

Since $\dim V$ is odd and $\dim \ker(T^2 + bT + cI)$ is even (by 5.33), we have that $\dim \text{im}(T^2 + bT + cI)$ is odd. Hence $\text{im}(T^2 + bT + cI)$ is a subspace of V that is invariant under T (by 7.11) and that has odd dimension less than $\dim V$. Our induction hypothesis now implies that T restricted to $\text{im}(T^2 + bT + cI)$ has an eigenvalue, which means that T has an eigenvalue. \square

§7.3 Upper-Triangular Matrices

Suppose $T \in \mathcal{L}(V)$. Recall that the matrix of T with respect to a basis $\{v_1, \dots, v_n\}$ of V is the $n \times n$ matrix whose entries a_{ij} are defined by

$$Tv_j = \sum_{i=1}^n a_{ij}v_i.$$

Notation. The notation $\mathcal{M}(T; \{v_1, \dots, v_n\})$ is used if the basis is not clear from the context.

Remark. The matrices of operators are square matrices.

The *diagonal* of a square matrix consists of the entries on the line from the upper left corner to the bottom right corner.

Definition 7.22 (Upper-triangular matrix). A square matrix is called *upper triangular* if all the entries below the diagonal are 0.

Typically we represent an upper-triangular matrix in the form

$$\begin{pmatrix} \lambda_1 & & * \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}$$

where the 0 indicates that all entries below the diagonal equal 0, and * denotes entries that we do not know or that are irrelevant to the questions being discussed.

The next result provides a useful connection between upper-triangular matrices and invariant subspaces.

Lemma 7.23 (Conditions for upper-triangular matrix). Suppose $T \in \mathcal{L}(V)$, $\{v_1, \dots, v_n\}$ is a basis of V . Then the following are equivalent:

- (i) The matrix of T with respect to $\{v_1, \dots, v_n\}$ is upper triangular.
- (ii) $\text{span}(v_1, \dots, v_k)$ is invariant under T for each $k = 1, \dots, n$.
- (iii) $Tv_k \in \text{span}(v_1, \dots, v_k)$ for each $k = 1, \dots, n$.

Proof.

(i) \implies (ii) Suppose $k \in \{1, \dots, n\}$. If $j \in \{1, \dots, n\}$, then

$$Tv_j \in \text{span}(v_1, \dots, v_j)$$

because the matrix of T with respect to $\{v_1, \dots, v_n\}$ is upper triangular. If $j \leq k$, then $\text{span}(v_1, \dots, v_j) \subset \text{span}(v_1, \dots, v_k)$, so

$$Tv_j \in \text{span}(v_1, \dots, v_k)$$

for each $j \in \{1, \dots, k\}$. Thus $\text{span}(v_1, \dots, v_k)$ is invariant under T .

(ii) \implies (iii) Suppose (ii) holds, so $\text{span}(v_1, \dots, v_k)$ is invariant under T for each $k = 1, \dots, n$. In particular, $Tv_k \in \text{span}(v_1, \dots, v_k)$ for each $k = 1, \dots, n$.

(iii) \implies (i) Suppose (iii) holds. Then when writing each Tv_k as a linear combination of basis vectors v_1, \dots, v_n , we need to use only v_1, \dots, v_k . Hence all entries under the diagonal of $\mathcal{M}(T)$ are 0, so $\mathcal{M}(T)$ is an upper-triangular matrix. \square

Lemma 7.24. *Suppose $T \in \mathcal{L}(V)$, V has a basis with respect to which T has an upper-triangular matrix with diagonal matrix with diagonal entries $\lambda_1, \dots, \lambda_n$. Then*

$$(T - \lambda_1 I) \cdots (T - \lambda_n I) = 0.$$

Proof. Let $\{v_1, \dots, v_n\}$ be a basis of V with respect to which T has an upper-triangular matrix with diagonal entries $\lambda_1, \dots, \lambda_n$.

- Considering the first column of $\mathcal{M}(T)$, we have

$$Tv_1 = \lambda_1 v_1 \implies (T - \lambda_1 I)v_1 = \mathbf{0},$$

which implies $(T - \lambda_1 I) \cdots (T - \lambda_m I)v_1 = \mathbf{0}$ for $m = 1, \dots, n$.

- Note that $(T - \lambda_2 I)v_2 \in \text{span}(v_1)$. Thus $(T - \lambda_1 I)(T - \lambda_2 I)v_2 = \mathbf{0}$ (by the previous paragraph), which implies $(T - \lambda_1 I) \cdots (T - \lambda_m I)v_2 = \mathbf{0}$ for $m = 2, \dots, n$.
- Note that $(T - \lambda_3 I)v_3 \in \text{span}(v_1, v_2)$. Thus by the previous paragraph, $(T - \lambda_1 I)(T - \lambda_2 I)(T - \lambda_3 I)v_3 = \mathbf{0}$, which implies $(T - \lambda_1 I) \cdots (T - \lambda_m I)v_3 = \mathbf{0}$ for $m = 3, \dots, n$.

Continuing this pattern, we see that

$$(T - \lambda_1 I) \cdots (T - \lambda_n I)v_k = \mathbf{0} \quad (k = 1, \dots, n).$$

Thus $(T - \lambda_1 I) \cdots (T - \lambda_n I)$ is the 0 operator because it is 0 on each vector in a basis of V . \square

Proposition 7.25. *Suppose $T \in \mathcal{L}(V)$ has an upper-triangular matrix with respect to some basis of V . Then the eigenvalues of T are precisely the entries on the diagonal of that upper-triangular matrix.*

Proof. Let $\{v_1, \dots, v_n\}$ be a basis of V with respect to which T has an upper-triangular matrix

$$\mathcal{M}(T) = \begin{pmatrix} \lambda_1 & & * \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}.$$

Since $Tv_1 = \lambda_1 v_1$, λ_1 is an eigenvalue of T . Suppose $k \in \{2, \dots, n\}$, then $(T - \lambda_k I)v_k \in \text{span}(v_1, \dots, v_{k-1})$, so $T - \lambda_k I$ maps $\text{span}(v_1, \dots, v_k)$ into $\text{span}(v_1, \dots, v_{k-1})$. Since

\square

The following result gives a necessary and sufficient condition to have an upper-triangular matrix.

Lemma 7.26. *Suppose V is finite-dimensional, $T \in \mathcal{L}(V)$. Then T has an upper-triangular matrix with respect to some basis of V if and only if the minimal polynomial equals $(z - \lambda_1) \cdots (z - \lambda_m)$ for some $\lambda_i \in \mathbf{F}$.*

Theorem 7.27. *Suppose V is finite-dimensional complex vector space, $T \in \mathcal{L}(V)$. Then T has an upper-triangular matrix with respect to some basis of V .*

Proof. The desired result follows immediately from 5.44 and the second version of the fundamental theorem of algebra (see 4.13). \square

§7.4 Diagonalisable Operators

Diagonal Matrices

Definition 7.28 (Diagonal matrix). A *diagonal matrix* is a square matrix that is 0 everywhere except possibly on the diagonal.

Remark. The entries on the diagonal are precisely the eigenvalues of the operator.

Definition 7.29 (Diagonalisable). An operator on V is called *diagonalisable* if the operator has a diagonal matrix with respect to some basis of V .

Remark. Diagonalisation may require a different basis.

Definition 7.30 (Eigenspace). Suppose $T \in \mathcal{L}(V)$, $\lambda \in \mathbf{F}$. The *eigenspace* of T corresponding to λ is the subspace of V defined by

$$E(\lambda, T) := \ker(T - \lambda I) = \{v \in V \mid Tv = \lambda v\}.$$

Remark. Hence $E(\lambda, T)$ is the set of all eigenvectors of T corresponding to λ , along with the $\mathbf{0}$ vector.

Proposition 7.31. Suppose $T \in \mathcal{L}(V)$, $\lambda_1, \dots, \lambda_m$ are distinct eigenvalues of T . Then

$$E(\lambda_1, T) + \dots + E(\lambda_m, T)$$

is a direct sum. Furthermore, if V is finite-dimensional, then

$$\dim E(\lambda_1, T) + \dots + \dim E(\lambda_m, T) \leq \dim V.$$

Conditions for Diagonalisability

Lemma 7.32 (Conditions equivalent to diagonalisability). *Suppose V is finite-dimensional, $T \in \mathcal{L}(V)$, $\lambda_1, \dots, \lambda_m$ are distinct eigenvalues of T . Then the following are equivalent:*

- (i) T is diagonalisable.
- (ii) V has a basis consisting of eigenvectors of T .
- (iii) $V = E(\lambda_1, T) \oplus \dots \oplus E(\lambda_m, T)$.
- (iv) $\dim V = \dim E(\lambda_1, T) + \dots + \dim E(\lambda_m, T)$.

Corollary 7.33. *Suppose V is finite-dimensional, $T \in \mathcal{L}(V)$ has $\dim V$ distinct eigenvalues. Then T is diagonalisable.*

Theorem 7.34. *Suppose V is finite-dimensional, $T \in \mathcal{L}(V)$. Then T is diagonalisable if and only if the minimal polynomial of T equals $(z - \lambda_1) \cdots (z - \lambda_m)$ for distinct $\lambda_1, \dots, \lambda_m \in \mathbf{F}$.*

Corollary 7.35. *Suppose $T \in \mathcal{L}(V)$ is diagonalisable, $U \leq V$ is invariant under T . Then $T|_U$ is a diagonalisable operator on U .*

Gershgorin Disk Theorem

Definition 7.36 (Gershgorin disks). Suppose $T \in \mathcal{L}(V)$, $\{v_1, \dots, v_n\}$ is a basis of V . Let A denote the matrix of T with respect to this basis. A **Gershgorin disk** of T with respect to the basis $\{v_1, \dots, v_n\}$ is a set of the form

$$\left\{ z \in \mathbf{F} \mid |z - a_{ii}| \leq \sum_{j=1, j \neq i}^n |a_{ij}| \right\},$$

where $i = 1, \dots, n$.

Theorem 7.37 (Gershgorin disk theorem). Suppose $T \in \mathcal{L}(V)$, $\{v_1, \dots, v_n\}$ is a basis of V . Then each eigenvalue of T is contained in some Gershgorin disk of T with respect to the basis $\{v_1, \dots, v_n\}$.

§7.5 Commuting Operators

Definition 7.38 (Commute). Two operators S and T on the same vector space *commute* if $ST = TS$. Two square matrices A and B of the same size commute if $AB = BA$.

Lemma 7.39 (commuting operators correspond to commuting matrices). Suppose $S, T \in \mathcal{L}(V)$ and $\{v_1, \dots, v_n\}$ is a basis of V . Then S and T commute if and only if $\mathcal{M}(S; \{v_1, \dots, v_n\})$ and $\mathcal{M}(T; \{v_1, \dots, v_n\})$ commute.

Lemma 7.40 (Eigenspace is invariant under commuting operator). Suppose $S, T \in \mathcal{L}(V)$ commute, $\lambda \in \mathbf{F}$. Then $E(\lambda, S)$ is invariant under T .

Proposition 7.41. Two diagonalizable operators on the same vector space have diagonal matrices with respect to the same basis if and only if the two operators commute.

Lemma 7.42 (Common eigenvector for commuting operators). Every pairs of commuting operators on a finite-dimensional nonzero complex vector space has a common eigenvector.

Lemma 7.43 (Commuting operators are simultaneously upper triangular). Suppose V is a finite-dimensional complex vector space, S and T are commuting operators on V . Then there is a basis of V with respect to which both S and T have upper-triangular matrices.

Proposition 7.44 (Eigenvalues of sum and product of commuting operators). Suppose V is a finite-dimensional complex vector space, S and T are commuting operators on V . Then

- (i) every eigenvalue of $S + T$ is an eigenvalue of S plus an eigenvalue of T ;
- (ii) every eigenvalue of ST is an eigenvalue of S times an eigenvalue of T .

Exercises

Exercise 7.1 ([Ax124] 5A Q1). Suppose $T \in \mathcal{L}(V)$, $U \leq V$. Prove that

- (i) if $U \subset \ker T$, then U is invariant under T ;
- (ii) if $\operatorname{im} T \subset U$, then U is invariant under T .

Solution.

- (i)
- (ii) Let $u \in U$. Then $Tu \in \operatorname{im} T \subset U$ so $Tu \in U$.

□

Exercise 7.2 ([Ax124] 5A Q4).

Exercise 7.3 ([Ax124] 5A Q8).

Exercise 7.4 ([Ax124] 5A Q11).

Exercise 7.5 ([Ax124] 5A Q13).

Exercise 7.6 ([Ax124] 5A Q28).

Exercise 7.7 ([Ax124] 5A Q32).

5B 2 7 10 11 13 17 18 22

8 Inner Product Spaces

§8.1 Inner Products and Norms

Inner Products

Definition 8.1. An *inner product* on V is a map $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbf{F}$ such that for all $u, v, w \in V$, $\lambda \in \mathbf{F}$,

(i) $\langle v, v \rangle \geq 0$, where equality holds if and only if $v = \mathbf{0}$; (positive definite)

(ii) $\langle u + v, w \rangle = \langle u, w \rangle + \langle v, w \rangle$ (sesquilinear)

$\langle \lambda u, v \rangle = \lambda \langle u, v \rangle$;

(iii) $\langle u, v \rangle = \overline{\langle v, u \rangle}$. (conjugate symmetric)

An *inner product space* $(V, \langle \cdot, \cdot \rangle)$ is a vector space V along with an inner product on V .

Notation. If the inner product on V is clear from context, we omit it and simply denote the inner product space as V .

Remark. Every real number equals its complex conjugate. Thus if we are dealing with a real vector space, then in (iii) we can dispense with the complex conjugate, so $\langle u, v \rangle = \langle v, u \rangle$ for all $u, v \in V$.

Example 8.2.

- The *Euclidean inner product* on \mathbf{F}^n is defined by

$$\langle (w_1, \dots, w_n), (z_1, \dots, z_n) \rangle = w_1 \overline{z_1} + \dots + w_n \overline{z_n}$$

for all $(w_1, \dots, w_n), (z_1, \dots, z_n) \in \mathbf{F}^n$.

- An inner product can be defined on the vector space $\mathcal{C}([-1, 1], \mathbb{R})$ by

$$\langle f, g \rangle = \int_{-1}^1 fg$$

for all $f, g \in \mathcal{C}([-1, 1], \mathbb{R})$.

Lemma 8.3 (Basic properties of inner product).

(i) For each fixed $u \in V$, the function that sends $v \mapsto \langle u, v \rangle$ is a linear map from V to \mathbf{F} .

(ii) $\langle 0, v \rangle = 0$ for every $v \in V$.

(iii) $\langle v, 0 \rangle = 0$ for every $v \in V$.

(iv) $\langle u, v + w \rangle = \langle u, v \rangle + \langle u, w \rangle$ for all $u, v, w \in V$.

(v) $\langle u, \lambda v \rangle = \bar{\lambda} \langle u, v \rangle$ for all $\lambda \in \mathbf{F}$, $u, v \in V$.

Proof.

(i) For $v \in V$, the linearity of $u \mapsto \langle u, v \rangle$ follows from the sesquilinearity of the inner product.

(ii) Every linear map takes $\mathbf{0}$ to 0 . Thus (ii) follows from (i).

(iii) If $v \in V$, by conjugate symmetry and (ii),

$$\langle v, 0 \rangle = \overline{\langle 0, v \rangle} = \bar{0} = 0.$$

(iv) Suppose $u, v, w \in V$. Then

$$\begin{aligned} \langle u, v + w \rangle &= \overline{\langle v + w, u \rangle} \\ &= \overline{\langle v, u \rangle + \langle w, u \rangle} \\ &= \overline{\langle v, u \rangle} + \overline{\langle w, u \rangle} \\ &= \langle u, v \rangle + \langle u, w \rangle. \end{aligned}$$

(v) Suppose $\lambda \in \mathbf{F}$, $u, v \in V$. Then

$$\begin{aligned} \langle u, \lambda v \rangle &= \overline{\langle \lambda v, u \rangle} \\ &= \overline{\lambda \langle v, u \rangle} \\ &= \bar{\lambda} \overline{\langle v, u \rangle} \\ &= \bar{\lambda} \langle u, v \rangle. \end{aligned}$$

□

Norms

Each inner product determines a norm.

Definition 8.4 (Norm). For $v \in V$, the *norm* of v is

$$\|v\| := \sqrt{\langle v, v \rangle}.$$

Lemma 8.5 (Basic properties of norm). Suppose $v \in V$.

- (i) $\|v\| = 0$ if and only if $v = \mathbf{0}$.
- (ii) $\|\lambda v\| = |\lambda| \|v\|$ for all $\lambda \in \mathbf{F}$.

Proof.

- (i) By positive definiteness of the inner product, $\langle v, v \rangle = 0$ if and only if $v = \mathbf{0}$. Take square root to get $\|v\| = 0$.
- (ii) Suppose $\lambda \in \mathbf{F}$. Then

$$\begin{aligned} \|\lambda v\|^2 &= \langle \lambda v, \lambda v \rangle \\ &= \lambda \langle v, \lambda v \rangle \\ &= \lambda \bar{\lambda} \langle v, v \rangle \\ &= |\lambda|^2 \|v\|^2. \end{aligned}$$

Taking square roots yields the desired equality.

□

Remark. Working with norms squared is usually easier than working directly with norms.

Now we come to a crucial definition.

Definition 8.6. $u, v \in V$ are *orthogonal* if $\langle u, v \rangle = 0$.

Lemma 8.7 (Orthogonality and $\mathbf{0}$).

- (i) $\mathbf{0}$ is orthogonal to every vector in V .
- (ii) $\mathbf{0}$ is the only vector in V that is orthogonal to itself.

Proof.

- (i) Recall that $\langle \mathbf{0}, v \rangle = 0$ for every $v \in V$.
- (ii) If $v \in V$ and $\langle v, v \rangle = 0$, then $v = \mathbf{0}$, by positive definiteness.

□

Lemma 8.8 (Pythagoras' theorem). *Suppose $u, v \in V$. If u and v are orthogonal, then*

$$\|u + v\|^2 = \|u\|^2 + \|v\|^2. \quad (8.1)$$

Proof. Suppose $\langle u, v \rangle = 0$. Then

$$\begin{aligned} \|u + v\|^2 &= \langle u + v, u + v \rangle \\ &= \langle u, u + v \rangle + \langle v, u + v \rangle \\ &= \langle u, u \rangle + \langle u, v \rangle + \langle v, u \rangle + \langle v, v \rangle \\ &= \|u\|^2 + 0 + \bar{0} + \|v\|^2 \\ &= \|u\|^2 + \|v\|^2 \end{aligned}$$

as desired. □

We now introduce a process known as *orthogonal decomposition*. Suppose $u, v \in V$, $u \neq \mathbf{0}$. Then the *orthogonal projection* of v onto u is

$$\text{proj}_u(v) := \frac{\langle v, u \rangle}{\langle u, u \rangle} u, \quad (8.2)$$

which is parallel to u . We check that $v - \text{proj}_u(v)$ and u are orthogonal:

$$\begin{aligned} \langle v - \text{proj}_u(v), u \rangle &= \langle v, u \rangle - \left\langle \frac{\langle v, u \rangle}{\langle u, u \rangle} u, u \right\rangle \\ &= \langle v, u \rangle - \frac{\langle v, u \rangle}{\langle u, u \rangle} \langle u, u \rangle = 0. \end{aligned}$$

Lemma 8.9 (Cauchy–Schwarz inequality). *Suppose $u, v \in V$. Then*

$$|\langle u, v \rangle| \leq \|u\| \|v\|, \quad (8.3)$$

where equality holds if and only if $u = \lambda v$ for some scalar λ .

Proof. If $u = \mathbf{0}$, then both sides of the desired inequality equal 0. Thus assume $u \neq \mathbf{0}$. Consider the orthogonal decomposition of v :

$$v = (v - \text{proj}_u(v)) + \text{proj}_u(v).$$

By the Pythagoras' theorem,

$$\|v\|^2 = \underbrace{\|v - \text{proj}_u(v)\|^2}_{\geq 0} + \|\text{proj}_u(v)\|^2,$$

so

$$\|v\| \geq \|\text{proj}_u(v)\| = \left| \frac{\langle v, u \rangle}{\langle u, u \rangle} \right| \|u\| = \frac{|\langle v, u \rangle|}{\|u\|}$$

and rearranging gives the desired inequality. Equality holds if and only if $v = \text{proj}_u(v)$, i.e.,

$$\frac{\langle v, u \rangle}{\langle u, u \rangle} u = v.$$

□

Lemma 8.10 (Triangle inequality). *Suppose $u, v \in V$. Then*

$$\|u + v\| \leq \|u\| + \|v\|, \quad (8.4)$$

where equality holds if and only if $u = \lambda v$ for some $\lambda \in \mathbb{R}_{\geq 0}$.

Proof. We have

$$\begin{aligned} \|u + v\|^2 &= \langle u + v, u + v \rangle \\ &= \langle u, u \rangle + \langle v, v \rangle + \langle u, v \rangle + \langle v, u \rangle \\ &= \langle u, u \rangle + \langle v, v \rangle + \langle u, v \rangle + \overline{\langle u, v \rangle} \\ &= \|u\|^2 + \|v\|^2 + 2 \operatorname{Re} \langle u, v \rangle \\ &\leq \|u\|^2 + \|v\|^2 + 2 |\langle u, v \rangle| \quad (1) \\ &\leq \|u\|^2 + \|v\|^2 + 2 \|u\| \|v\| \quad (2) \\ &= (\|u\| + \|v\|)^2, \end{aligned}$$

where (2) follows from the Cauchy–Schwarz inequality. Taking square roots of both sides of the above inequality gives the desired inequality.

Equality holds if and only if equality holds in (1) and (2), i.e.,

$$\langle u, v \rangle = \|u\| \|v\|.$$

If $u = \lambda v$ for $\lambda \in \mathbb{R}_{\geq 0}$, then the above equation holds. Conversely, suppose the above equation holds. Then equality in the Cauchy–Schwarz inequality implies that $u = \lambda v$ for some scalar λ . By the above equation, λ must be a non-negative real number, completing the proof. \square

Corollary 8.11 (Reverse triangle inequality).

Lemma 8.12 (Parallelogram equality). *Suppose $u, v \in V$. Then*

$$\|u + v\|^2 + \|u - v\|^2 = 2 (\|u\|^2 + \|v\|^2). \quad (8.5)$$

Proof. We have

$$\begin{aligned} \|u + v\|^2 + \|u - v\|^2 &= \langle u + v, u + v \rangle + \langle u - v, u - v \rangle \\ &= (\|u\|^2 + \|v\|^2 + \langle u, v \rangle + \langle v, u \rangle) + (\|u\|^2 + \|v\|^2 - \langle u, v \rangle - \langle v, u \rangle) \\ &= 2 (\|u\|^2 + \|v\|^2) \end{aligned}$$

as desired. \square

§8.2 Orthonormal Bases

Orthonormal Bases

Definition 8.13 (Orthonormal basis). $\{e_1, \dots, e_n\} \subset V \setminus \{0\}$ is *orthonormal* if

- (i) $\|e_i\| = 1$;
- (ii) the vectors are pairwise orthogonal.

If additionally $\{e_1, \dots, e_n\}$ is a basis of V , then $\{e_1, \dots, e_n\}$ is a *orthonormal basis* of V .

Lemma 8.14. Suppose $\{e_1, \dots, e_n\}$ is a orthonormal set of vectors in V . Then

$$\|a_1e_1 + \dots + a_n e_n\|^2 = |a_1|^2 + \dots + |a_n|^2$$

for all $a_1, \dots, a_n \in \mathbf{F}$.

Proof. By the Pythagoras' theorem,

$$\begin{aligned} \|a_1e_1 + \dots + a_n e_n\|^2 &= \|a_1e_1\|^2 + \dots + \|a_n e_n\|^2 \\ &= |a_1|^2 \|e_1\|^2 + \dots + |a_n|^2 \|e_n\|^2 \\ &= |a_1|^2 + \dots + |a_n|^2 \end{aligned}$$

since each $\|e_i\| = 1$. □

The result above has the following important corollary.

Corollary 8.15. Every orthonormal set of vectors is linearly independent.

Proof. Suppose $\{e_1, \dots, e_n\}$ is an orthonormal set of vectors in V . Suppose $a_1, \dots, a_n \in \mathbf{F}$ are such that

$$a_1e_1 + \dots + a_n e_n = \mathbf{0}.$$

By the previous result,

$$|a_1|^2 + \dots + |a_n|^2 = 0,$$

so $a_1 = \dots = a_n = 0$. Hence e_1, \dots, e_n are linearly independent. □

Hence every orthonormal set of vectors in V of length $\dim V$ is an orthonormal basis of V .

Now we come to an important inequality.

Lemma 8.16 (Bessel's inequality). Suppose $\{e_1, \dots, e_n\}$ is an orthonormal set of vectors in V . If $v \in V$ then

$$|\langle v, e_1 \rangle|^2 + \dots + |\langle v, e_n \rangle|^2 \leq \|v\|^2. \quad (8.6)$$

Proof. Let $v \in V$. For $i = 1, \dots, n$, consider the orthogonal projection of v onto e_i :

$$\begin{aligned} v &= (v - \text{proj}_{e_i}(v)) + \text{proj}_{e_i}(v) \\ &= \left(v - \frac{\langle v, e_i \rangle}{\langle e_i, e_i \rangle} e_i \right) + \frac{\langle v, e_i \rangle}{\langle e_i, e_i \rangle} e_i \\ &= (v - \langle v, e_i \rangle e_i) + \langle v, e_i \rangle e_i. \end{aligned}$$

Then by Pythagoras' theorem,

$$\begin{aligned} \|v\|^2 &= \|v - \langle v, e_i \rangle e_i\|^2 + \|\langle v, e_i \rangle e_i\|^2 \\ &= \|v - \langle v, e_i \rangle e_i\|^2 + |\langle v, e_i \rangle|^2. \end{aligned}$$

Write

$$\begin{aligned} v &= \text{proj}_{e_1}(v) + \dots + \text{proj}_{e_n}(v) + w \\ &= \langle v, e_1 \rangle e_1 + \dots + \langle v, e_n \rangle e_n + w \end{aligned}$$

for some $w \in V$. Note that for $i = 1, \dots, n$,

$$\begin{aligned} \langle v, e_i \rangle &= \langle \langle v, e_i \rangle e_i + \langle w, e_i \rangle, e_i \rangle \\ &= \langle v, e_i \rangle + \langle w, e_i \rangle \end{aligned}$$

which implies $\langle w, e_i \rangle = 0$, so w is orthogonal to e_1, \dots, e_n . Thus e_1, \dots, e_n, w are pairwise orthogonal. By Pythagoras' theorem,

$$\begin{aligned} \|v\|^2 &= |\langle v, e_1 \rangle|^2 + \dots + |\langle v, e_n \rangle|^2 + \underbrace{\|w\|^2}_{\geq 0} \\ &\geq |\langle v, e_1 \rangle|^2 + \dots + |\langle v, e_n \rangle|^2 \end{aligned}$$

as desired. Equality holds for orthonormal bases (as we will see later). \square

The next result states that a vector can be expressed as a linear combination of an orthonormal basis. Usually we write $v = \sum_{i=1}^n a_i v_i$, but with orthonormal basis we can just take $a_i = \langle v, e_i \rangle$.

Lemma 8.17. *Suppose $\{e_1, \dots, e_n\}$ is an orthonormal basis of V , and $u, v \in V$. Then*

(i) $v = \langle v, e_1 \rangle e_1 + \dots + \langle v, e_n \rangle e_n$;

(ii) $\|v\|^2 = |\langle v, e_1 \rangle|^2 + \dots + |\langle v, e_n \rangle|^2$; (Parseval's identity)

(iii) $\langle u, v \rangle = \langle u, e_1 \rangle \overline{\langle v, e_1 \rangle} + \dots + \langle u, e_n \rangle \overline{\langle v, e_n \rangle}$.

Proof.

(i) Since e_1, \dots, e_n is a basis of V , there exist $a_1, \dots, a_n \in \mathbf{F}$ such that

$$v = a_1 e_1 + \dots + a_n e_n.$$

Since e_1, \dots, e_n are orthonormal, taking the inner product of both sides with e_i gives

$$\langle v, e_i \rangle = a_i \quad (i = 1, \dots, n).$$

Hence we are done.

(ii) Apply Pythagoras' theorem to (i).

(iii) Taking the inner product of u with each side of (i),

$$\begin{aligned} \langle u, v \rangle &= \langle u, \langle v, e_1 \rangle e_1 + \dots + \langle v, e_n \rangle e_n \rangle \\ &= \langle u, \langle v, e_1 \rangle e_1 \rangle + \dots + \langle u, \langle v, e_n \rangle e_n \rangle \\ &= \overline{\langle v, e_1 \rangle} \langle u, e_1 \rangle + \dots + \overline{\langle v, e_n \rangle} \langle u, e_n \rangle. \end{aligned}$$

□

Gram–Schmidt Procedure

The *Gram–Schmidt procedure* is a method for constructing orthonormal basis, by turning a linearly independent set into an orthonormal set with the same span as the original set. It guarantees the existence of orthonormal bases.

Theorem 8.18 (Gram–Schmidt procedure). *Suppose v_1, \dots, v_n are linearly independent in V . Define*

$$u_i = \begin{cases} v_1 & (i = 1) \\ v_i - \text{proj}_{u_1}(v_i) - \dots - \text{proj}_{u_{i-1}}(v_i) & (i = 2, \dots, n) \end{cases}$$

and let

$$e_i = \frac{u_i}{\|u_i\|}.$$

Then $\{e_1, \dots, e_n\}$ is an orthonormal set of vectors in V such that

$$\text{span}(v_1, \dots, v_i) = \text{span}(e_1, \dots, e_i)$$

for $i = 1, \dots, n$.

Proof. Induct on i . For $i = 1$, since $e_1 = \frac{u_1}{\|u_1\|}$ we have $\|e_1\| = 1$, and $\text{span}(v_1) = \text{span}(e_1)$ because e_1 is a non-zero multiple of v_1 .

Suppose the desired result holds for $i - 1$; that is, the set $\{e_1, \dots, e_{i-1}\}$ generated by the above procedure is an orthonormal set, and

$$\text{span}(v_1, \dots, v_{i-1}) = \text{span}(e_1, \dots, e_{i-1}).$$

Since v_1, \dots, v_n are linearly independent, we have $v_i \notin \text{span}(v_1, \dots, v_{i-1})$. Thus $v_i \notin \text{span}(e_1, \dots, e_{i-1}) = \text{span}(u_1, \dots, u_{i-1})$, which implies that $u_i \neq \mathbf{0}$ (so we are not dividing by 0); thus $\|e_i\| = 1$.

We now check that e_1, \dots, e_i is an orthonormal set. For $j = 1, \dots, i - 1$,

$$\begin{aligned} \langle e_i, e_j \rangle &= \left\langle \frac{u_i}{\|u_i\|}, \frac{u_j}{\|u_j\|} \right\rangle \\ &= \frac{1}{\|u_i\| \|u_j\|} \langle u_i, u_j \rangle \\ &= \frac{1}{\|u_i\| \|u_j\|} \left\langle v_i - \text{proj}_{u_1}(v_i) - \dots - \text{proj}_{u_j}(v_i) - \dots - \text{proj}_{u_{i-1}}(v_i), u_j \right\rangle \\ &= \frac{1}{\|u_i\| \|u_j\|} \left\langle v_i - \frac{\langle v_i, u_1 \rangle}{\langle u_1, u_1 \rangle} u_1 - \dots - \frac{\langle v_i, u_j \rangle}{\langle u_j, u_j \rangle} u_j - \dots - \frac{\langle v_i, u_{i-1} \rangle}{\langle u_{i-1}, u_{i-1} \rangle} u_{i-1}, u_j \right\rangle \\ &= \frac{1}{\|u_i\| \|u_j\|} \left(\langle v_i, u_j \rangle - \left\langle \frac{\langle v_i, u_j \rangle}{\langle u_j, u_j \rangle} u_j, u_j \right\rangle \right) \\ &= \frac{1}{\|u_i\| \|u_j\|} \left(\langle v_i, u_j \rangle - \langle v_i, u_j \rangle \right) = 0 \end{aligned}$$

so e_i is orthogonal to e_1, \dots, e_{i-1} .

□

Corollary 8.19. *Every finite-dimensional inner product space has an orthonormal basis.*

Corollary 8.20. *Suppose V is finite-dimensional. Then every orthonormal set of vectors in V can be extended to an orthonormal basis of V .*

Proof.

□

Corollary 8.21. *Suppose V is finite-dimensional. Then every orthonormal list of vectors in V can be extended to an orthonormal basis of V .*

Proposition 8.22. *Suppose V is finite-dimensional, $T \in \mathcal{L}(V)$. Then T has an upper-triangular matrix with respect to some orthonormal basis of V if and only if the minimal polynomial of T equals $(z - \lambda_1) \cdots (z - \lambda_n)$ for some $\lambda_i \in \mathbf{F}$.*

Theorem 8.23 (Schur's theorem). *Every operator on a finite-dimensional complex inner product space has an upper-triangular matrix with respect to some orthonormal basis.*

Linear Functionals on Inner Product Spaces

Theorem 8.24 (Riesz representation theorem). *Suppose V is finite-dimensional and ϕ is a linear functional on V . Then there exists a unique vector $v \in V$ such that*

$$\phi(u) = \langle u, v \rangle$$

for every $u \in V$.

§8.3 Orthogonal Complements and Minimisation Problems

Orthogonal Complements

Definition 8.25. The *orthogonal complement* of $U \subset V$ is

$$U^\perp := \{v \in V \mid \langle u, v \rangle = 0, \forall u \in U\}.$$

Lemma 8.26 (Properties of orthogonal complement).

- (i) If $U \subset V$, then $U^\perp \leq V$.
- (ii) $\{\mathbf{0}\}^\perp = V$, $V^\perp = \{\mathbf{0}\}$.
- (iii) If $U \subset V$, then $U \cap U^\perp \subset \{\mathbf{0}\}$.
- (iv) If $G \subset H \subset V$, then $H^\perp \subset G^\perp$.

Proof.

(i) Let $U \subset V$. Use the subspace test (4.10):

(i) $\langle u, \mathbf{0} \rangle = 0$ for every $u \in U$, so $\mathbf{0} \in U^\perp$.

(ii) Let $v, w \in U^\perp$. For every $u \in U$,

$$\langle u, v + w \rangle = \langle u, v \rangle + \langle u, w \rangle = 0 + 0 = 0 \implies v + w \in U^\perp$$

so U^\perp is closed under addition.

(iii) Let $v \in U^\perp$, $\lambda \in \mathbf{F}$. For every $u \in U$,

$$\langle u, \lambda v \rangle = \bar{\lambda} \langle u, v \rangle = \bar{\lambda} \cdot 0 = 0 \implies \lambda v \in U^\perp$$

so U^\perp is closed under scalar multiplication.

(ii)

$$v \in V \iff \langle \mathbf{0}, v \rangle = 0 \iff v \in \{\mathbf{0}\}^\perp$$

$$v \in V^\perp \iff \langle v, v \rangle = 0 \iff v = \mathbf{0}$$

(iii)

(iv)

□

Lemma 8.27. Suppose $U \leq V$ is finite-dimensional. Then

$$V = U \oplus U^\perp$$

and thus $\dim U^\perp = \dim V - \dim U$. In addition,

$$U = (U^\perp)^\perp.$$

Lemma 8.28. Suppose $U \leq V$ is finite-dimensional. Then

$$U^\perp = \{\mathbf{0}\} \iff U = V.$$

Definition 8.29 (Orthogonal projection). Suppose $U \leq V$ is finite-dimensional. The *orthogonal projection* is the operator $P_U \in \mathcal{L}(V)$ defined as follows: For each $v \in V$, write $v = u + w$ for some $u \in U, w \in U^\perp$. Then let $P_U v = u$.

Remark. Suppose $u \in V$ with

Lemma 8.30 (Properties of orthogonal projection). Suppose $U \leq V$ is finite-dimensional.

- (i) $P_U \in \mathcal{L}(V)$.
- (ii) $P_U u = u$ for every $u \in U$, $P_U w = \mathbf{0}$ for every $w \in U^\perp$.
- (iii) $\text{im } P_U = U$, $\ker P_U = U^\perp$.
- (iv) $v - P_U v \in U^\perp$ for every $v \in V$.
- (v) $P_U^2 = P_U$.
- (vi) $\|P_U v\| \leq \|v\|$ for every $v \in V$.
- (vii) If e_1, \dots, e_n is an orthonormal basis of U , and $v \in V$, then

$$P_U v = \langle v, e_1 \rangle e_1 + \dots + \langle v, e_n \rangle e_n.$$

(i) to (v) hold for projection maps, (vi) and (vii) are specific to orthogonal projections.

Minimisation Problems

Given a subspace U of V and a point $v \in V$, we want to find a point $u \in U$ that minimises $\|v - u\|$. The next result shows that $u = P_U v$ is the unique solution of this minimisation problem.

Theorem 8.31 (Minimising distance to a subspace). *Suppose $U \leq V$ is finite-dimensional. Fix $v \in V$. Then for all $u \in U$*

$$\|v - P_U v\| \leq \|v - u\|, \quad (8.7)$$

where equality holds if and only if $u = P_U v$.

Proof. We have

$$\begin{aligned} \|v - P_U v\|^2 &\leq \|v - P_U v\|^2 + \|P_U v - u\|^2 \quad [\because \|P_U v - u\|^2 \geq 0] \\ &= \end{aligned}$$

□

Pseudoinverse

Restriction of a linear map to obtain a bijective map.

Lemma 8.32. *Suppose V is finite-dimensional, and $T \in \mathcal{L}(V)$. Then $T|_{(\ker T)^\perp}$ is an injective map from $(\ker T)^\perp$ to $\text{im } T$.*

Definition 8.33 (Pseudoinverse). *Suppose V is finite-dimensional, and $T \in \mathcal{L}(V)$. The **pseudoinverse** $T^+ \in \mathcal{L}(W, V)$ of T is the linear map from W to V defined by*

$$T^+w = \left(T|_{(\ker T)^\perp}\right)^{-1} P_{\text{im } T}w$$

for each $w \in W$.

The pseudoinverse is also called the Moore–Penrose inverse.

The pseudoinverse behaves much like an inverse, as we will see.

Lemma 8.34 (Properties of pseudoinverse). *Suppose V is finite-dimensional, and $T \in \mathcal{L}(V)$.*

- (i) *If T is invertible, then $T^+ = T^{-1}$.*
- (ii) *$TT^+ = P_{\text{im } T}$.*
- (iii) *$T^+T = P_{(\ker T)^\perp}$.*

Theorem 8.35 (Pseudoinverse provides best approximate solution or best solution). *Suppose V is finite-dimensional, $T \in \mathcal{L}(V, W)$, and $w \in W$.*

(i) *If $v \in V$, then*

$$\left\|T(T^+w) - w\right\| \leq \|Tv - w\|, \quad (8.8)$$

where equality holds if and only if $v \in T^+w + \ker T$.

(ii) *If $v \in T^+w + \ker T$, then*

$$\left\|T^+w\right\| \leq \|v\|, \quad (8.9)$$

where equality holds if and only if $v = T^+w$.

Exercises

Exercise 8.1 ([Ax124] 6A Q1). Show that if $v_1, \dots, v_m \in V$, then

$$\sum_{j=1}^m \sum_{k=1}^m \langle v_j, v_k \rangle \geq 0.$$

Solution.

□

Exercise 8.2 ([Ax124] 6A Q2). Suppose $S \in \mathcal{L}(V)$. Define $\langle \cdot, \cdot \rangle_1$ by

$$\langle u, v \rangle_1 = \langle Su, Sv \rangle$$

for all $u, v \in V$. Show that $\langle \cdot, \cdot \rangle_1$ is an inner product on V if and only if S is injective.

Exercise 8.3 ([Ax124] 6A Q4). Suppose $T \in \mathcal{L}(V)$ is such that $\|Tv\| \leq \|v\|$ for every $v \in V$. Prove that $T - \sqrt{2}I$ is injective.

Exercise 8.4 ([Ax124] 6A Q5).

Exercise 8.5 ([Ax124] 6A Q6).

Exercise 8.6 ([Ax124] 6A Q9).

Exercise 8.7 ([Ax124] 6A Q14).

Exercise 8.8 ([Ax124] 6A Q20).

Exercise 8.9 ([Ax124] 6A Q26).

Exercise 8.10 ([Ax124] 6A Q27).

Exercise 8.11 ([Ax124] 6B Q1).

Exercise 8.12 ([Ax124] 6B Q3).

Exercise 8.13 ([Ax124] 6B Q5).

Exercise 8.14 ([Ax124] 6B Q6(a)).

Exercise 8.15 ([Ax124] 6B Q9).

Exercise 8.16 ([Ax124] 6B Q10).

Exercise 8.17 ([Ax124] 6B Q13).

Exercise 8.18 ([Ax124] 6B Q14).

Exercise 8.19 ([Ax124] 6B Q15).

Exercise 8.20 ([Ax124] 6B Q16).

Exercise 8.21 ([Ax124] 6B Q17).

Exercise 8.22 ([Ax124] 6B Q23).

9 Operators on Inner Product Spaces

§9.1 Self-Adjoint and Normal Operators

Adjoint

Definition 9.1. Suppose $T \in \mathcal{L}(V, W)$. The *adjoint* is the function $T^* : V \rightarrow V$ such that

$$\langle Tv, w \rangle = \langle v, T^*w \rangle \quad (v \in V, w \in W).$$

We check that if $T \in \mathcal{L}(V, W)$, then $T^* \in \mathcal{L}(V, W)$; that is, the adjoint of a linear map is a linear map.

Lemma 9.2 (Properties of adjoint). Suppose $T \in \mathcal{L}(V, W)$. Then

- (i) $(S + T)^* = S^* + T^*$ for all $S \in \mathcal{L}(V, W)$
- (ii) $(\lambda T)^* = \bar{\lambda}T^*$ for all $\lambda \in \mathbf{F}$
- (iii) $(T^*)^* = T$
- (iv) $(ST)^* = T^*S^*$ for all $S \in \mathcal{L}(W, U)$, where U is a finite-dimensional inner product space over \mathbf{F}
- (v) $I^* = I$, where I is the identity operator on V
- (vi) if T is invertible, then T^* is invertible, and $(T^*)^{-1} = (T^{-1})^*$

Lemma 9.3 (Kernel and image of T^*). Suppose $T \in \mathcal{L}(V, W)$. Then

- (i) $\ker T^* = (\operatorname{im} T)^\perp$
- (ii) $\operatorname{im} T^* = (\ker T)^\perp$
- (iii) $\ker T = (\operatorname{im} T^*)^\perp$
- (iv) $\operatorname{im} T = (\ker T^*)^\perp$

Definition 9.4 (Conjugate transpose).

Proposition 9.5.

Self-Adjoint Operators

Definition 9.6 (Self-adjoint operator). An operator $T \in \mathcal{L}(V)$ is called *self-adjoint* if $T = T^*$.

That is, an operator $T \in \mathcal{L}(V)$ is self-adjoint if and only if

$$\langle Tv, w \rangle = \langle v, Tw \rangle$$

for all $v, w \in V$.

Lemma 9.7. *Every eigenvalue of a self-adjoint operator is real.*

Lemma 9.8. *Suppose V is a complex inner product space, and $T \in \mathcal{L}(V)$. Then*

$$\langle Tv, v \rangle = 0 \forall v \in V \iff T = 0.$$

Remark. do not hold for real inner product spaces

Lemma 9.9. *Suppose V is a complex inner product space, and $T \in \mathcal{L}(V)$. Then*

$$T \text{ is self-adjoint} \iff \langle Tv, v \rangle \in \mathbb{R} \forall v \in V.$$

Remark. do not hold for real inner product spaces

Lemma 9.10. *Suppose T is a self-adjoint operator on V . Then*

$$\langle Tv, v \rangle = 0 \forall v \in V \iff T = 0.$$

Normal Operators

Definition 9.11 (Normal operator). An operator on an inner product space is *normal* if it commutes with its adjoint.

That is, $T \in \mathcal{L}(V)$ is normal if $TT^* = T^*T$.

Remark. Every self-adjoint operator is normal, but not vice versa.

Lemma 9.12 (Characterisation of normal operators). Suppose $T \in \mathcal{L}(V)$. Then

$$T \text{ is normal} \iff \|Tv\| = \|T^*v\| \quad \forall v \in V.$$

Lemma 9.13. Suppose $T \in \mathcal{L}(V)$. Then

- (i) $\ker T = \ker T^*$
- (ii) $\operatorname{im} T = \operatorname{im} T^*$
- (iii) $V = \ker T \oplus \operatorname{im} T$
- (iv) $T - \lambda I$ is normal for every $\lambda \in \mathbf{F}$
- (v) if $v \in V$ and $\lambda \in \mathbf{F}$, then $Tv = \lambda v$ if and only if $T^*v = \bar{\lambda}v$

Proposition 9.14. Suppose $T \in \mathcal{L}(V)$ is normal. Then the eigenvalues of T corresponding to distinct eigenvalues are orthogonal.

Proposition 9.15. Suppose $\mathbf{F} = \mathbb{C}$ and $T \in \mathcal{L}(V)$. Then T is normal if and only if there exist commuting self-adjoint operators A and B such that $T = A + iB$.

§9.2 Spectral Theorem

Real Spectral Theorem

Complex Spectral Theorem

§9.3 Positive Operators

§9.4 Isometries, Unitary Operators, and Matrix Factorisation

Isometries

Unitary Operators

QR Factorisation

Cholesky Factorisation

§9.5 Singular Value Decomposition

Singular Values

SVD for Linear Maps and for Matrices

§9.6 Consequences of Singular Value Decomposition

Norms of Linear Maps

Approximation by Linear Maps with Lower-Dimensional Range

Polar Decomposition

Operators Applied to Ellipsoids and Parallelepipeds

Volume via Singular Values

Properties of an Operator as Determined by Its Eigenvalues

IV

Real Analysis

Real analysis deals with the real numbers and real-valued functions of a real variable.

A great part of analysis deals with inequalities and error terms. This is evident from the very beginning, in the theory of epsilons and deltas. Instead of obtaining precise values, it is sufficient to show that epsilon and delta are within a certain range. In order to show convergence, we just need to show that the error terms are small. Thus, there is often no perfect bound or best approximation, and there need not be; all that is needed is for the bound or the approximation to be good enough.

10 Real and Complex Number Systems

Summary

- Supremum, infimum.
- Construction and properties of the real field \mathbb{R} .
- Construction and properties of the complex field \mathbb{C} .
- Construction and properties of the Euclidean space \mathbb{R}^n .

§10.1 Ordered Sets and Boundedness

Definitions

Let S be a set.

Definition 10.1. An *order* on S is a binary relation $<$ such that

- (i) for all $x, y \in S$, exactly one of $x < y$, $x = y$, or $y < x$ holds; (trichotomy)
- (ii) if $x, y, z \in S$ are such that $x < y$ and $y < z$, then $x < z$. (transitivity)

S is an *ordered set* if it has an order; denote it by $(S, <)$.

Notation. We write $x \leq y$ if $x < y$ or $x = y$. We define $>$ and \geq in the obvious way.

Definition 10.2 (Boundedness). Let $E \subset S$, where S is an ordered set.

- (i) E is *bounded above* if there exists $\beta \in S$ such that $x \leq \beta$ for all $x \in E$; we call β an *upper bound* of E .
- (ii) E is *bounded below* if there exists $\beta \in S$ such that $x \geq \beta$ for all $x \in E$; we call β a *lower bound* of E .

E is *bounded* in S if it is bounded above and below.

Definition 10.3 (Supremum, infimum). We say $\alpha \in S$ is the *supremum* of E if

- (i) α is an upper bound for E ;
- (ii) if $\beta < \alpha$ then β is not an upper bound of E , i.e. $\exists x \in S$ s.t. $x > \beta$ (least upper bound).

Likewise, we say $\alpha \in S$ is the *infimum* of E if

- (i) α is a lower bound for E ;
- (ii) if $\beta > \alpha$ then β is not a lower bound of E , i.e. $\exists x \in S$ s.t. $x < \beta$ (greatest lower bound).

Remark. It is not necessary for the supremum and infimum of E to be in E .

Lemma 10.4 (Uniqueness of supremum). *If E has a supremum, then it is unique.*

Proof. Suppose α and β be suprema of E .

Since β is a supremum, it is an upper bound for E . Since α is a supremum, then it is the *least* upper bound, so $\alpha \leq \beta$. Interchanging the roles of α and β gives $\beta \leq \alpha$. Hence $\alpha = \beta$. \square

Since the supremum and infimum are unique, we can give them a notation.

Notation. Denote the supremum of E by $\sup E$, the infimum by $\inf E$.

Example 10.5. Let $E = \left\{ \frac{1}{n} \mid n \in \mathbb{N} \right\}$. Then $\sup E = 1$, $\inf E = 0$.

Proof. It is clear that 1 is an upper bound for E . Suppose $\beta < 1$. Since $1 \in E$, evidently β is not an upper bound for E . Hence $\sup E = 1$.

It is clear that 0 is a lower bound for E . Suppose $\beta > 0$. Pick $n = \left\lfloor \frac{1}{\beta} \right\rfloor + 1$, then $\beta > \frac{1}{n}$, so β is not a lower bound for E . Hence $\inf E = 0$. \square

Least-upper-bound Property

Definition 10.6. An ordered set S has the *least-upper-bound property* (l.u.b.) if every non-empty subset of S that is bounded above has a supremum in S .

We define the *greatest-lower-bound property* similarly.

Proposition 10.7. *Suppose S is an ordered set. If S has the least-upper-bound property, then S has the greatest-lower-bound property.*

Proof. Suppose S has the least-upper-bound property. Let non-empty $B \subset S$ be bounded below. We want to show that $\inf B \in S$.

Let $L \subset S$ be the set of all lower bounds of B ; that is,

$$L = \{y \in S \mid y \leq x \forall x \in B\}.$$

Since B is bounded below, B has a lower bound, so $L \neq \emptyset$. Since every $x \in B$ is an upper bound of L , L is bounded above. By the least-upper-bound property of S , we have that $\sup L \in S$.

Claim. $\inf B = \sup L$.

To show that $\sup L = \inf B$ (greatest lower bound), we need to show that (i) $\sup L$ is a lower bound of B , (ii) and $\sup L$ is the greatest of the lower bounds.

- (i) Suppose $\gamma < \sup L$, then γ is not an upper bound of L . Since B is the set of upper bounds of L , $\gamma \notin B$. Considering the contrapositive, if $\gamma \in B$, then $\gamma \geq \sup L$. Hence $\sup L$ is a lower bound of B , and thus $\sup L \in L$.
- (ii) If $\sup L < \beta$ then $\beta \notin L$, since $\sup L$ is an upper bound of L . In other words, $\sup L$ is a lower bound of B , but β is not if $\beta > \sup L$. This means that $\sup L$ is the greatest of the lower bounds.

Hence $\inf B = \sup L \in S$. □

Corollary 10.8. *If S has the greatest-lower-bound property, then it has the least-upper-bound property. Hence S has the least-upper-bound property if and only if S has the greatest-lower-bound property.*

Properties of Suprema and Infima

This section discusses some fundamental properties of the supremum that will be useful in this text. There is a corresponding set of properties of the infimum that the reader should formulate for himself.

The next result shows that a set with a supremum contains numbers arbitrarily close to its supremum.

Lemma 10.9 (Approximation property). *Let $S \subset \mathbb{R}$ be non-empty, $b = \sup S$. Then for every $a < b$ there exists $x \in S$ such that*

$$a < x \leq b.$$

Proof. We first show $x \leq b$. Since $b = \sup S$ is an upper bound of S , $x \leq b$ for all $x \in S$.

We now show there exist $x \in S$ such that $a < x$. Suppose otherwise, for a contradiction, that $x \leq a$ for every $x \in S$. Then a would be an upper bound for S . But since $a < b$ and b is the supremum, this means a is smaller than the least upper bound, a contradiction. \square

For the rest of this section, suppose S has the least-upper-bound property.

Lemma 10.10 (Additive property). *Given non-empty $A, B \subset S$, let*

$$C = \{x + y \mid x \in A, y \in B\}.$$

If each of A and B has a supremum, then C has a supremum, and

$$\sup C = \sup A + \sup B.$$

Proof. Let $a = \sup A$, $b = \sup B$. Let $z \in C$, then $z = x + y$ for some $x \in A$, $y \in B$. Then

$$z = x + y \leq a + b,$$

so $a + b$ is an upper bound for C . Since C is non-empty and bounded above, by the lub property of S , C has a supremum in S .

Let $c = \sup C$. To show that $a + b = c$, we need to show that (i) $a + b \geq c$, and (ii) $a + b \leq c$.

(i) Since c is the *least* upper bound for C , and $a + b$ is an upper bound for C , we must have that $c \leq a + b$.

(ii) Choose any $\varepsilon > 0$. By 10.9 there exist $x \in A$ and $y \in B$ such that

$$a - \varepsilon < x, \quad b - \varepsilon < y.$$

Adding these inequalities gives

$$a + b - 2\varepsilon < x + y \leq c.$$

Thus $a + b < c + 2\varepsilon$ for every $\varepsilon > 0$. Hence $a + b \leq c$.

\square

Lemma 10.11 (Comparison property). *Let non-empty $A, B \subset S$ such that $a \leq b$ for every $a \in A$,*

$b \in B$. If B has a supremum, then A has a supremum, and

$$\sup A \leq \sup B.$$

Proof. Let $\beta = \sup B$. Since β is a supremum for B , then $b \leq \beta$ for all $b \in B$.

Let $a \in A$ and choose any $b \in B$. Since $a \leq b$ and $b \leq \beta$, $a \leq \beta$. Thus β is an upper bound for A .

Since A is non-empty and bounded above, by the lub property of S , A has a supremum in S ; let $\alpha = \sup A$. Since β is an upper bound for A , and α is the *least* upper bound for A , we have that $\alpha \leq \beta$, as desired. \square

Lemma 10.12. Let $B \subset S$ be non-empty and bounded below. Let

$$A = -B := \{-b \mid b \in B\}.$$

Then A is non-empty and bounded above. Furthermore, $\inf B$ exists, and $\inf B = -\sup A$.

Proof. Since B is non-empty, so is A . Since B is bounded below, let β be a lower bound for B . Then $b \geq \beta$ for all $b \in B$, which implies $-b \leq -\beta$ for all $b \in B$. Hence $a \leq -\beta$ for all $a \in A$, so $-\beta$ is an upper bound for A .

Since A is non-empty and bounded above, by the lub property of S , A has a supremum. Then $a \leq \sup A$ for all $a \in A$, so $b \geq -\sup A$ for all $b \in B$. Thus $-\sup A$ is a lower bound for B .

Also, we saw before that if β is a lower bound for B then $-\beta$ is an upper bound for A . Then $-\beta \geq \sup A$ (since $\sup A$ is the least upper bound), so $\beta \leq -\sup A$. Therefore $-\sup A$ is the greatest lower bound of B . \square

Ordered Fields

Definition 10.13 (Ordered field). A field F is an **ordered field** if there exists an order $<$ on F such that for all $x, y, z \in F$,

- (i) if $y < z$ then $x + y < x + z$;
- (ii) if $x > 0$ and $y > 0$ then $xy > 0$.

If $x > 0$, we call x *positive*; if $x < 0$, x is *negative*.

All the familiar rules for working with inequalities apply in every ordered field: Multiplication by positive [negative] quantities preserves [reverses] inequalities, no square is negative, etc. The following result lists some of these.

Lemma 10.14 (Basic properties). *Let F be an ordered field, $x, y, z \in F$.*

- (i) *If $x > 0$ then $-x < 0$, and vice versa.*
- (ii) *If $x > 0$ and $y < z$ then $xy < xz$.*
- (iii) *If $x < 0$ and $y < z$ then $xy > xz$.*
- (iv) *If $x \neq 0$ then $x^2 > 0$. In particular, $1 > 0$.*
- (v) *If $0 < x < y$ then $0 < \frac{1}{y} < \frac{1}{x}$.*

Proof.

(i) If $x > 0$ then $0 = -x + x > -x + 0$, so that $-x < 0$.

If $x < 0$ then $0 = -x + x < -x + 0$, so that $-x > 0$.

(ii) Since $z > y$, we have $z - y > y - y = 0$, so $x(z - y) > 0$. Hence

$$xz = x(z - y) + xy > 0 + xy = xy.$$

(iii) By (i) and (ii),

$$-[x(z - y)] = (-x)(z - y) > 0,$$

so that $x(z - y) < 0$. Hence $xz < xy$.

(iv) If $x > 0$, part (ii) of the above definition gives $x^2 > 0$.

If $x < 0$, then $-x > 0$ so $(-x)^2 > 0$. But $x^2 = (-x)^2$.

Since $1 = 1^2$, $1 > 0$.

(v) If $y > 0$ and $v \leq 0$, then $yv \leq 0$. But $y\left(\frac{1}{y}\right) = 1 > 0$, so $\frac{1}{y} > 0$. Likewise, $\frac{1}{x} > 0$.

Multiplying both sides of the inequality $x < y$ by the positive quantity $\left(\frac{1}{x}\right)\left(\frac{1}{y}\right)$, we obtain $\frac{1}{y} < \frac{1}{x}$.

□

§10.2 Real Numbers

Problems with \mathbb{Q}

\mathbb{Q} has some problems, the first of which being *algebraic incompleteness*: there exists equations with coefficients in \mathbb{Q} but do not have solutions in \mathbb{Q} (in fact \mathbb{R} has this problem too, but \mathbb{C} is algebraically complete, by the fundamental theorem of algebra).

Lemma 10.15. $x^2 - 2 = 0$ has no solution in \mathbb{Q} .

Proof. Suppose, for a contradiction, that $x^2 - 2 = 0$ has a solution $x = \frac{p}{q}$, $q \neq 0$. We also assume $\frac{p}{q}$ is in lowest terms; that is, p, q are coprime. Squaring both sides gives $\frac{p^2}{q^2} = 2$, or $p^2 = 2q^2$. Observe that p^2 is even, so p is even; let $p = 2m$ for some integer m . Then this implies $4m^2 = 2q^2$, or $2m^2 = q^2$. Similarly, q^2 is even so q is even.

Since p and q share a common factor of 2, we have reached a contradiction. \square

The second problem is *analytic incompleteness*: there exists a sequence of rational numbers that approach a point that is not in \mathbb{Q} ; for example, the sequence

$$1, 1.4, 1.41, 1.414, 1.4142, \dots$$

tends to the the irrational number $\sqrt{2}$.

Continuing from the above lemma,

Lemma 10.16. Let

$$A = \{p \in \mathbb{Q} \mid p > 0, p^2 < 2\},$$

$$B = \{p \in \mathbb{Q} \mid p > 0, p^2 > 2\}.$$

Then A contains no largest number, and B contains no smallest number.

Proof. Prove by construction. We associate with each rational $p > 0$ the number

$$q = p - \frac{p^2 - 2}{p + 2} = \frac{2p + 2}{p + 2}$$

and so

$$q^2 - 2 = \frac{2(p^2 - 2)}{(p + 2)^2}.$$

For any $p \in A$, $q > p$ and $q \in A$ since $q^2 < 2$, so A has no largest number.

For any $p \in B$, $q < p$ and $q \in B$ since $q^2 > 2$, so B has no smallest number. \square

Proposition 10.17. \mathbb{Q} does not have the least-upper-bound property.

Proof. In the previous result, note that B is the set of all upper bounds of A , and B does not have a smallest element. Hence $A \subset \mathbb{Q}$ is bounded above but A has no least upper bound in \mathbb{Q} . \square

Real Field

The sole objective of this subsection is to prove the following result.

Theorem 10.18 (Existence of real field). *There exists an ordered field \mathbb{R} that*

- (i) *contains \mathbb{Q} as a subfield, and*
- (ii) *has the least-upper-bound property (also known as the completeness axiom).*

We want to construct \mathbb{R} from \mathbb{Q} ; one method to do so is using Dedekind cuts.

Definition 10.19 (Dedekind cut). $\alpha \subset \mathbb{Q}$ is a *Dedekind cut*, if

- (i) $\alpha \neq \emptyset, \alpha \neq \mathbb{Q}$; (non-trivial)
- (ii) if $p \in \alpha, q \in \mathbb{Q}$ and $q < p$, then $q \in \alpha$;
- (iii) if $p \in \alpha$, then $p < r$ for some $r \in \alpha$.

Remark. Note that (iii) simply says that α has no largest member; (ii) implies two facts which will be used freely:

- If $p \in \alpha$ and $q \notin \alpha$, then $p < q$.
- If $r \notin \alpha$ and $r < s$, then $s \notin \alpha$.

Example 10.20. Let $r \in \mathbb{Q}$ and define

$$\alpha_r := \{p \in \mathbb{Q} \mid p < r\}.$$

We now check that this is indeed a Dedekind cut.

- (i) $p = 1 + r \notin \alpha_r$ thus $\alpha_r \neq \mathbb{Q}$. $p = r - 1 \in \alpha_r$ thus $\alpha_r \neq \emptyset$.
- (ii) Suppose that $q \in \alpha_r$ and $q' < q$. Then $q' < q < r$ which implies that $q' < r$ thus $q' \in \alpha_r$.
- (iii) Suppose that $q \in \alpha_r$. Consider $\frac{q+r}{2} \in \mathbb{Q}$ and $q < \frac{q+r}{2} < r$. Thus $\frac{q+r}{2} \in \alpha_r$.

This example shows that every rational r corresponds to a Dedekind cut α_r .

Example 10.21. $\sqrt[3]{2}$ is not rational, but it is real. $\sqrt[3]{2}$ corresponds to the cut

$$\alpha = \{p \in \mathbb{Q} \mid p^3 < 2\}.$$

- (i) Trivial.
- (ii) If $q < p$, by the monotonicity of the cubic function, this implies that $q^3 < p^3 < 2$ thus $q \in \alpha$.
- (iii) If $p \in \alpha$, consider $\left(p + \frac{1}{n}\right)^3 < 2$.

Definition 10.22. The set of real numbers, denoted by \mathbb{R} , is the set of all Dedekind cuts:

$$\mathbb{R} := \{\alpha \subset \mathbb{Q} \mid \alpha \text{ is a Dedekind cut}\}.$$

Proposition 10.23. \mathbb{R} has an order, where $\alpha < \beta$ is defined to mean that $\alpha \subsetneq \beta$.

Proof. Simply check if this is a valid order (by checking for trichotomy and transitivity). \square

Proposition 10.24. The ordered set \mathbb{R} has the least-upper-bound property.

Proof. Let non-empty $A \subset \mathbb{R}$ be bounded above. Let $\beta \in \mathbb{R}$ be an upper bound of A . We want to show that A has a supremum in \mathbb{R} .

Let

$$\gamma = \bigcup_{\alpha \in A} \alpha.$$

Then $p \in \gamma$ if and only if $p \in \alpha$ for some $\alpha \in A$.

Claim. $\gamma \in \mathbb{R}$ and $\gamma = \sup A$.

We first prove that $\gamma \in \mathbb{R}$ by checking that it is a Dedekind cut:

- (i) Since $A \neq \emptyset$, there exists $\alpha_0 \in A$. Since $\alpha_0 \in \mathbb{R}$, it is a Dedekind cut so $\alpha_0 \neq \emptyset$. Since $\alpha_0 \subset \gamma$, $\gamma \neq \emptyset$.
Since $\alpha \subset \beta$ for every $\alpha \in A$, the union of $\alpha \in A$ must be a subset of β ; thus $\gamma \subset \beta$. Hence $\gamma \neq \mathbb{Q}$.
- (ii) Let $p \in \gamma$. Then $p \in \alpha_1$ for some $\alpha_1 \in A$. If $q < p$, then $q \in \alpha_1$ (since α_1 is a Dedekind cut). Hence $q \in \gamma$.
- (iii) If $r \in \alpha_1$ is so chosen that $r > p$, we see that $r \in \gamma$ (since $\alpha_1 \subset \gamma$).

Next we prove that $\gamma = \sup A$, by checking that (i) γ is an upper bound of A , (ii) γ is the *least* of the upper bounds.

- (i) It is clear that $\alpha \leq \gamma$ for every $\alpha \in A$.
- (ii) Suppose $\delta < \gamma$. Then there exists $s \in \gamma$ such that $s \notin \delta$. Since $s \in \gamma$, $s \in \alpha$ for some $\alpha \in A$. Hence $\delta < \alpha$, so δ is not an upper bound of A .

\square

Remark. The l.u.b. property of \mathbb{R} is also known as the *completeness axiom* of \mathbb{R} .

We now define operations on \mathbb{R} .

Definition 10.25 (Addition). Given $\alpha, \beta \in \mathbb{R}$, define addition as

$$\alpha + \beta := \{r \in \mathbb{Q} \mid r = a + b, a \in \alpha, b \in \beta\}.$$

We first check if the above definition makes sense. We want to show that addition on \mathbb{R} is closed: for all $\alpha, \beta \in \mathbb{R}$, $\alpha + \beta \in \mathbb{R}$.

Proof. We check that $\alpha + \beta$ is a Dedekind cut:

(i) Since $\alpha \neq \emptyset$ and $\beta \neq \emptyset$, there exists $a \in \alpha$ and $b \in \beta$. Hence $r = a + b \in \alpha + \beta$ so $\alpha + \beta \neq \emptyset$.

Since $\alpha \neq \mathbb{Q}$ and $\beta \neq \mathbb{Q}$, there exist $c \notin \alpha$ and $d \notin \beta$. Thus $r' = c + d > a + b$ for any $a \in \alpha, b \in \beta$, so $r' \notin \alpha + \beta$. Hence $\alpha + \beta \neq \mathbb{Q}$.

(ii) Suppose that $r \in \alpha + \beta$ and $r' < r$. We want to show that $r' \in \alpha + \beta$.

$r = a + b$ for some $a \in \alpha, b \in \beta$. Then $r' - a < b$. Since $\beta \in \mathbb{R}$, $r' - a \in \beta$ so $r' - a = b_1$ for some $b_1 \in \beta$. Hence $r' = a + b_1 \in \alpha + \beta$.

(iii) Suppose $r \in \alpha + \beta$, so $r = a + b$ for some $a \in \alpha, b \in \beta$. Since α, β are Dedekind cuts, there exist $a' \in \alpha, b' \in \beta$ with $a < a'$ and $b < b'$. Then $r = a + b < a' + b' \in \alpha + \beta$. We define $r' = a' + b' \in \alpha + \beta$ with $r < r'$.

□

Lemma 10.26.

(i) Addition on \mathbb{R} is commutative: $\alpha + \beta = \beta + \alpha$ for all $\alpha, \beta \in \mathbb{R}$.

(ii) Addition on \mathbb{R} is associative: $\alpha + (\beta + \gamma) = (\alpha + \beta) + \gamma$ for all $\alpha, \beta, \gamma \in \mathbb{R}$.

(iii) Additive identity: Define $0^* := \{p \in \mathbb{Q} \mid p < 0\}$. Then $\alpha + 0^* = \alpha$ for all $\alpha \in \mathbb{R}$.

(iv) Additive inverse: Fix $\alpha \in \mathbb{R}$, define $\beta = \{p \in \mathbb{Q} \mid \exists r > 0, -p - r \notin \alpha\}$. Then $\alpha + \beta = 0^*$.

Remark. Recall that to prove that two sets are equal, show double inclusion.

Proof.

(i) We need to show that $\alpha + \beta \subset \beta + \alpha$ and $\beta + \alpha \subset \alpha + \beta$.

Let $r \in \alpha + \beta$. Then $r = a + b$ for $a \in \alpha$ and $b \in \beta$. Thus $r = b + a$ since $+$ is commutative on \mathbb{Q} . Hence $r \in \beta + \alpha$. Therefore $\alpha + \beta \subset \beta + \alpha$.

Similarly, $\beta + \alpha \subset \alpha + \beta$.

Therefore $\alpha + \beta = \beta + \alpha$.

(ii) Let $r \in \alpha + (\beta + \gamma)$. Then $r = a + (b + c)$ where $a \in \alpha, b \in \beta, c \in \gamma$. Thus $r = (a + b) + c$ by associativity of $+$ on \mathbb{Q} . Therefore $r \in (\alpha + \beta) + \gamma$, hence $\alpha + (\beta + \gamma) \subset (\alpha + \beta) + \gamma$.

Similarly, $(\alpha + \beta) + \gamma \subset \alpha + (\beta + \gamma)$.

(iii) It is clear that 0^* is a Dedekind cut.

Let $r \in \alpha + 0^*$. Then $r = a + p$ for some $a \in \alpha, p \in 0^*$. Thus $r = a + p < a + 0 = a$ so $r \in \alpha$. Hence $\alpha + 0^* \subset \alpha$.

Let $r \in \alpha$. Then there exists $r' \in \alpha$ where $r' > r$. Thus $r - r' < 0$, so $r - r' \in 0^*$. We see that $r = r' + (r - r')$ where $r' \in \alpha, r - r' \in 0^*$. Hence $\alpha \subset \alpha + 0^*$.

(iv) Fix some $\alpha \in \mathbb{R}$. We first show that β is a Dedekind cut.

(i) Let $s \notin \alpha$, let $p = -s - 1$. Then $-p - 1 \notin \alpha$. Hence $p \in \beta$, so $\beta \neq \emptyset$.

Let $q \in \alpha$. Then $-q \notin \beta$ so $\beta \neq \mathbb{Q}$.

(ii) Let $p \in \beta$. Then there exists $r > 0$ such that $-p - r \notin \alpha$. If $q < p$, then $-q - r > -p - r$ so $-q - r \notin \alpha$. Hence $q \in \beta$.

(iii) Let $t = p + \frac{r}{2}$. Then $t > p$, and $-t - \frac{r}{2} = -p - r \notin \alpha$. Hence $t \in \beta$.

Let $r \in \alpha, s \in \beta$. Then $-s \notin \alpha$. This implies $r < -s$ (since α is closed downwards) so $r + s < 0$. Hence $\alpha + \beta \subset 0^*$.

To prove the opposite inclusion, let $v \in 0^*$, and let $w = -\frac{v}{2}$. Then $w > 0$. By the Archimedean property on \mathbb{Q} , there exists $n \in \mathbb{N}$ such that $nw \in \alpha$ but $(n + 1)w \notin \alpha$. Let $p = -(n + 2)w$. Then

$$-p - w = (n + 2)w - w = (n + 1)w \notin \alpha$$

so $p \in \beta$. Since $v = nw + p$ where $nw \in \alpha, p \in \beta, v \in \alpha + \beta$. Hence $0^* \subset \alpha + \beta$. □

Notation. β is denoted by the more familiar notation $-\alpha$.

Lemma 10.27. *If $\alpha, \beta, \gamma \in \mathbb{R}$ and $\beta < \gamma$, then $\alpha + \beta < \alpha + \gamma$.*

Proof. □

We say that a Dedekind cut α is *positive* if $0 \in \alpha$, and *negative* if $0 \notin \alpha$. If α is neither positive nor negative, then $\alpha = 0^*$.

Multiplication is a little more bothersome than addition in the present context, since products of negative rationals are positive. For this reason we confine ourselves first to \mathbb{R}^+ (the set of all $\alpha \in \mathbb{R}$ with $\alpha > 0^*$).

Definition 10.28. Given $\alpha, \beta \in \mathbb{R}^+$, define multiplication as

$$\alpha\beta := \{p \in \mathbb{Q} \mid p \leq rs, r \in \alpha, s \in \beta, r, s > 0\}.$$

We also define $1^* := \{q \in \mathbb{Q} \mid q < 1\}$.

As again, check if the above definition makes sense. We want to show that multiplication on \mathbb{R}^+ is closed: for all $\alpha, \beta \in \mathbb{R}, \alpha\beta \in \mathbb{R}$.

Proof. Check that $\alpha\beta$ is a Dedekind cut.

(i) $\alpha \neq \emptyset$ means there exists $r \in \alpha, r > 0$. Similarly, $\beta \neq \emptyset$ means there exists $s \in \beta, s > 0$. Then $rs \in \mathbb{Q}$ and $rs \leq rs$, so $rs \in \alpha\beta$. Hence $\alpha\beta \neq \emptyset$.

$\alpha \neq \mathbb{Q}$ means there exists $r' \notin \alpha$ such that $r' > r$ for all $r \in \alpha$. Similarly $\beta \neq \mathbb{Q}$ means there exists $s' \in \beta$ such that $s' > s$ for all $s \in \beta$. Then $r's' > rs$ for all $r \in \alpha, s \in \beta$, so $r's' \notin \alpha\beta$. Hence $\alpha\beta \neq \mathbb{Q}$.

(ii) Let $p \in \alpha\beta$. Then $p \leq ab$ for some $a \in \alpha, b \in \beta, a, b > 0$.

If $q < p$, then $q < p \leq ab$ so $q \in \alpha\beta$.

(iii) Let $p \in \alpha\beta$. Then $p \leq ab$ for some $a \in \alpha, b \in \beta, a, b > 0$. Pick $a' \in \alpha$ and $b' \in \beta$ with $a' > a$ and $b' > b$. Form $a'b' > ab \geq p, a'b' \leq a'b'$ means $a'b' \in \alpha \cdot \beta$. □

We now complete the definition of multiplication by setting $\alpha 0^* = 0^* = 0^* \alpha$, and by setting

$$\alpha\beta = \begin{cases} (-\alpha)(-\beta) & \alpha < 0^*, \beta < 0^*, \\ -[(-\alpha)\beta] & \alpha < 0^*, \beta > 0^*, \\ -[\alpha(-\beta)] & \alpha > 0^*, \beta < 0^*. \end{cases}$$

where we make negative numbers positive, multiply, and then negate them as needed.

Lemma 10.29.

- (i) *Multiplication on \mathbb{R} is commutative: $\alpha\beta = \beta\alpha$ for all $\alpha, \beta \in \mathbb{R}$.*
- (ii) *Multiplication on \mathbb{R} is associative: $(\alpha\beta)\gamma = \alpha(\beta\gamma)$ for all $\alpha, \beta, \gamma \in \mathbb{R}$.*
- (iii) *Multiplicative identity: $1\alpha = \alpha$ for all $\alpha \in \mathbb{R}$.*
- (iv) *Multiplicative inverse: If $\alpha \in \mathbb{R}$, $\alpha \neq 0^*$, then there exists $\beta \in \mathbb{R}$ such that $\alpha\beta = 1^*$.*

We associate each $r \in \mathbb{Q}$ with the set

$$r^* = \{p \in \mathbb{Q} \mid p < r\}.$$

It is obvious that each r^* is a cut; that is, $r^* \in \mathbb{R}$.

Proposition 10.30. *The replacement of $r \in \mathbb{Q}$ by the corresponding “rational cuts” $r^* \in \mathbb{R}$ preserves sums, products, and order. That is, for all $r^*, s^* \in \mathbb{R}$,*

- (i) $r^* + s^* = (r + s)^*$;
- (ii) $r^* s^* = (rs)^*$;
- (iii) $r^* < s^*$ if and only if $r < s$.

Proof.

- (i) Let $p \in r^* + s^*$. Then $p = u + v$ for some $u \in r^*$, $v \in s^*$, where $u < r$, $v < s$. Then $p < r + s$. Hence $p \in (r + s)^*$, so $r^* + s^* \subset (r + s)^*$.

Let $p \in (r + s)^*$. Then $p < r + s$. Let $t = \frac{(r+s)-p}{2}$, and let

$$r' = r - t, \quad s' = s - t.$$

Since $t > 0$, $r' < r$ so $r' \in r^*$; $s' < s$ so $s' \in s^*$. Then $p = r' + s'$, so $p \in r^* + s^*$. Hence $(r + s)^* \subset r^* + s^*$.

(ii)

- (iii) Suppose $r < s$. Then $r \in s^*$, but $r \notin r^*$. Hence $r^* < s^*$.

Conversely, suppose $r^* < s^*$. Then there exists $p \in s^*$ such that $p \in r^*$. Hence $r \leq p < s$, so $r < s$.

□

This shows that the ordered field \mathbb{Q} is isomorphic to the ordered field $\mathbb{Q}^* = \{q^* \mid q \in \mathbb{Q}\}$ whose elements are rational cuts. It is this identification of \mathbb{Q} with \mathbb{Q}^* which allows us to regard \mathbb{Q} as a subfield of \mathbb{R} .

Remark. In fact, \mathbb{R} is the only ordered field with the l.u.b. property. Hence any other ordered field with the l.u.b. property is isomorphic to \mathbb{R} .

Therefore we have proven [10.18](#).

Properties of \mathbb{R}

Proposition 10.31 (Archimedean property). *For any $x \in \mathbb{R}^+$, $y \in \mathbb{R}$, there exists $n \in \mathbb{N}$ such that*

$$nx > y.$$

Proof. Suppose, for a contradiction, that $nx \leq y$ for all $n \in \mathbb{N}$. Then y is an upper bound of the set

$$A = \{nx \mid n \in \mathbb{N}\}.$$

Since $A \subset \mathbb{R}$ is non-empty and bounded above, by the l.u.b. property of \mathbb{R} , A has a supremum in \mathbb{R} , say $\alpha = \sup A$.

Consider $\alpha - x$. Since $\alpha - x < \alpha = \sup A$, $\alpha - x$ is not an upper bound of A . Then $\alpha - x \leq n_0x$ for some $n_0 \in \mathbb{N}$; rearranging gives $\alpha \leq (n_0 + 1)x$. This implies that α is not an upper bound of A , which contradicts the fact that α is the supremum of A . \square

Corollary 10.32. *Let $\varepsilon > 0$. Then there exists $n \in \mathbb{N}$ such that $0 < \frac{1}{n} < \varepsilon$.*

Proof. In 10.31, take $x = \varepsilon$ and $y = 1$. \square

Proposition 10.33 (\mathbb{Q} is dense in \mathbb{R}). *For any $x, y \in \mathbb{R}$ with $x < y$, there exists $p \in \mathbb{Q}$ such that*

$$x < p < y.$$

Proof. We prove by construction (construct the required p from the given x and y).

Since $x < y$, we have $y - x > 0$. By 10.32, there exists $n \in \mathbb{N}$ such that

$$\frac{1}{n} < y - x.$$

Consider the set comprising multiples of $\frac{1}{n}$:

$$E = \left\{ \frac{k}{n} \mid k \in \mathbb{N} \right\}.$$

Since E is unbounded, choose the first multiple $m \in \mathbb{N}$ such that $\frac{m}{n} > x$.

Claim. $x < \frac{m}{n} < y$.

It suffices to show that $\frac{m}{n} < y$. If not, then

$$\frac{m-1}{n} < x \quad \text{and} \quad \frac{m}{n} > y,$$

where the first inequality follows from the minimality of m . But these two statements combined imply that $\frac{1}{n} > y - x$, a contradiction. \square

Proposition 10.34 (\mathbb{R} is closed under taking roots). *For every $x \in \mathbb{R}^+$ and every $n \in \mathbb{N}$, there exists a unique $y \in \mathbb{R}^+$ so that $y^n = x$.*

We call the number y the positive n -th root of x , and denote it by $\sqrt[n]{x}$ or $x^{\frac{1}{n}}$.

Proof. Let $x \in \mathbb{R}^+$, fix $n \in \mathbb{N}$.

Existence Let

$$E = \{t \in \mathbb{R}^+ \mid t^n < x\}.$$

Claim. $y = \sup E$ satisfies $y^n = x$.

We first show that E has a supremum.

(i) Let $t = \frac{x}{1+x}$. Then $0 \leq t < 1$, so $t^n \leq t < x$ implies $t^n < x$. Hence $t \in E$, so $E \neq \emptyset$.

(ii) We claim that $1+x$ is an upper bound for E .

If $t > 1+x$, then $t^n \geq t > x$ implies $t^n > x$, so $t \notin E$. [This is the contrapositive of $t \in E \implies t \leq 1+x$.] Hence $1+x$ is an upper bound of E , so E is bounded above.

Hence E has a supremum; let $y = \sup E$.

To prove that $y^n = x$, we show that both the inequalities $y^n < x$ and $y^n > x$ lead to a contradiction. Consider the identity $b^n - a^n = (b-a)(b^{n-1} + b^{n-2}a + \cdots + a^{n-1})$. If $0 < a < b$, then

$$b^n - a^n < (b-a)nb^{n-1}. \quad (1)$$

Case 1: $y^n < x$.

Idea. We can find a small $h > 0$ such that $(y+h)^n < x$.

Choose h so that $0 < h < 1$ and

$$h < \frac{x - y^n}{n(y+1)^{n-1}}.$$

In (1), take $b = y+h$, $a = y$. Then

$$\begin{aligned} (y+h)^n - y^n &< hn(y+h)^{n-1} \\ &< hn(y+1)^{n-1} \\ &< \frac{x - y^n}{n(y+1)^{n-1}} n(y+1)^{n-1} \\ &= x - y^n. \end{aligned}$$

Thus $(y+h)^n < x$, and $y+h \in E$. Since $y+h > y$, this contradicts the fact that y is an upper bound of E .

Case 2: $y^n > x$.

Idea. Similarly, we can find a small $k > 0$ such that $(y-k)^n > x$.

Let

$$k = \frac{y^n - x}{ny^{n-1}}.$$

Then $0 < k < y$, by (1). If $t \geq y - k$,

$$\begin{aligned} y^n - t^n &\leq y^n - (y - k)^n \\ &< kny^{n-1} \\ &= \frac{y^n - x}{ny^{n-1}} \\ &= y^n - x. \end{aligned}$$

Thus $t^n > x$, and $t \notin E$. It follows that $y - k$ is an upper bound of E . But $y - k < y$, which contradicts the fact that y is the *least* upper bound of E .

Uniqueness Suppose, for a contradiction, that there exist distinct y_1, y_2 which are both n -th roots of x . WLOG assume that $0 < y_1 < y_2$. Then taking the n -th power gives $y_1^n < y_2^n$.

Since y_1 is a n -th root of x , then $x = y_1^n$, so $x < y_2^n$ implies $x \neq y_2^n$. Hence y_2 cannot be a n -th root of x , a contradiction. \square

Corollary 10.35. *If $a, b \in \mathbb{R}^+$ and $n \in \mathbb{N}$, then*

$$(ab)^{\frac{1}{n}} = a^{\frac{1}{n}} b^{\frac{1}{n}}.$$

Proof. Let $\alpha = a^{\frac{1}{n}}, \beta = b^{\frac{1}{n}}$. Then

$$ab = \alpha^n \beta^n = (\alpha\beta)^n$$

since multiplication is commutative. The uniqueness assertion of the previous result shows that

$$(ab)^{\frac{1}{n}} = \alpha\beta = a^{\frac{1}{n}} b^{\frac{1}{n}}.$$

\square

Lemma 10.36. *If $x \in \mathbb{R}^+$ and $m, n \in \mathbb{N}$, then*

$$\left(x^{\frac{1}{n}}\right)^m = \left(x^m\right)^{\frac{1}{n}}.$$

Proof. Exercise. \square

We can now define *rational exponents* x^r , where $x > 0$ and $r \in \mathbb{Q}$.

Definition 10.37 (Rational exponents). For $x > 0$ and $m, n \in \mathbb{N}$, define

$$x^{\frac{m}{n}} := \left(x^{\frac{1}{n}}\right)^m \quad \text{and} \quad x^{-\frac{m}{n}} := \frac{1}{x^{\frac{m}{n}}}.$$

(We also define $x^0 = 1$.)

We need to check that the above definition of x^r is well defined. That is, if $m, n, p, q \in \mathbb{N}$ are such that $\frac{m}{n} = \frac{p}{q}$, then $\left(x^{\frac{1}{n}}\right)^m = \left(x^{\frac{1}{q}}\right)^p$. To see this, note that $mq = np$ and

$$\left(\left(x^{\frac{1}{n}}\right)^m\right)^q = \left(x^{\frac{1}{n}}\right)^{mq} = \left(x^{\frac{1}{n}}\right)^{np} = x^p.$$

Thus $(x^{\frac{1}{n}})^m$ is the q -th root of x^p , i.e.,

$$(x^{\frac{1}{n}})^m = (x^p)^{\frac{1}{q}}.$$

Lemma 10.38 (Properties of rational exponents).

- (i) If $a > 0$ and $r, s \in \mathbb{Q}$, then $a^{r+s} = a^r a^s$ and $(a^r)^s = a^{rs}$.
- (ii) If $0 < a < b$ and $r \in \mathbb{Q}$ with $r > 0$, then $a^r < b^r$.
- (iii) If $a > 1$, $r, s \in \mathbb{Q}$ with $r < s$, then $a^r < a^s$.

The next result shows that real numbers can be approximated to any desired degree of accuracy by rational numbers with finite decimal representations.

Proposition 10.39. Let $x \geq 0$. Then for every integer $n \geq 1$ there exists a finite decimal $r_n = a_0.a_1a_2 \cdots a_n$ such that

$$r_n \leq x < r_n + \frac{1}{10^n}.$$

Proof. We prove by construction (construct the required finite decimal from x).

Let

$$S = \{k \in \mathbb{Z} \mid k \leq x\}.$$

S is non-empty (since $0 \in S$), and S is bounded above by x . Hence by the lub property of \mathbb{R} , S has a supremum in \mathbb{R} , say $a_0 = \sup S$. It is easily verified that $a_0 \in S$, so a_0 is a non-negative integer. We call a_0 the *greatest integer in x* , and write $a_0 = \lfloor x \rfloor$. Clearly we have

$$a_0 \leq x < a_0 + 1.$$

Now let $a_1 = \lfloor 10(x - a_0) \rfloor$. Since $0 \leq 10(x - a_0) < 10$, we have $0 \leq a_1 \leq 9$ and

$$a_1 \leq 10x - 10a_0 < a_1 + 1.$$

In other words, a_1 is the largest integer satisfying the inequalities

$$a_0 + \frac{a_1}{10} \leq x < a_0 + \frac{a_1 + 1}{10}.$$

More generally, having chosen a_1, \dots, a_{n-1} with $0 \leq a_i \leq 9$, let a_n be the largest integer satisfying the inequalities

$$a_0 + \frac{a_1}{10} + \cdots + \frac{a_n}{10^n} \leq x < a_0 + \frac{a_1}{10} + \cdots + \frac{a_n + 1}{10^n}.$$

Then $0 \leq a_n \leq 9$ and we have

$$r_n \leq x < r_n + \frac{1}{10^n},$$

where $r_n = a_0.a_1a_2 \cdots a_n$. □

Furthermore, it is easy to verify that $x = \sup_{n \in \mathbb{N}} r_n$.

Extended Real Number System

Definition 10.40 (Extended real number system). The *extended real number system* is defined to be the union

$$\overline{\mathbb{R}} := \mathbb{R} \cup \{-\infty, +\infty\},$$

where we preserve the original order in \mathbb{R} , and define $-\infty < x < +\infty$ for all $x \in \mathbb{R}$.

Defining $\overline{\mathbb{R}}$ is convenient since the following result holds.

Proposition 10.41. Any non-empty $E \subset \overline{\mathbb{R}}$ has a supremum and infimum in $\overline{\mathbb{R}}$.

Proof. If E is bounded above in \mathbb{R} , then by the l.u.b. property of \mathbb{R} , it has a supremum in $\mathbb{R} \subset \overline{\mathbb{R}}$. If E is not bounded above in \mathbb{R} , then $\sup E = +\infty \in \overline{\mathbb{R}}$.

Exactly the same remarks apply to lower bounds. □

$\overline{\mathbb{R}}$ does not form a field, but it is customary to make the following conventions for arithmetic on $\overline{\mathbb{R}}$:

(i) If $x \in \mathbb{R}$ then

$$x + \infty = +\infty, \quad x - \infty = -\infty, \quad \frac{x}{+\infty} = \frac{x}{-\infty} = 0.$$

(ii) If $x > 0$ then

$$x \cdot (+\infty) = +\infty, \quad x \cdot (-\infty) = -\infty.$$

If $x < 0$ then

$$x \cdot (+\infty) = -\infty, \quad x \cdot (-\infty) = +\infty.$$

When it is desired to make the distinction between real numbers on the one hand and the symbols $+\infty$ and $-\infty$ on the other quite explicit, the former are called *finite*.

§10.3 Complex Field

Lemma 10.42. Let $(a, b), (c, d) \in \mathbb{R}^2$. Define addition and multiplication on \mathbb{R}^2 as

$$\begin{aligned}(a, b) + (c, d) &= (a + c, b + d), \\ (a, b)(c, d) &= (ac - bd, ad + bc).\end{aligned}$$

Then \mathbb{R}^2 is a field, with additive identity $(0, 0)$ and multiplicative identity $(1, 0)$.

We call this structure \mathbb{C} , the **complex field**; its elements are called *complex numbers*.

Proof. Check the field axioms. □

The next result shows that the complex numbers of the form $(a, 0)$ have the same arithmetic properties as the corresponding real numbers a . We can therefore identify $(a, 0) \in \mathbb{C}$ with $a \in \mathbb{R}$. This identification implies that \mathbb{R} is a subfield of \mathbb{C} .

Lemma 10.43. For any $a, b \in \mathbb{R}$, we have

$$\begin{aligned}(a, 0) + (b, 0) &= (a + b, 0), \\ (a, 0)(b, 0) &= (ab, 0).\end{aligned}$$

Proof. Exercise. □

You may have noticed that we have defined the complex numbers without referring to the mysterious square root of -1 . We now show that the notation (a, b) is equivalent to the more customary $a + bi$.

Define the imaginary number $i := (0, 1)$. See that

$$i^2 = (0, 1)(0, 1) = (-1, 0) = -1.$$

Lemma 10.44. For $a, b \in \mathbb{R}$, $(a, b) = a + bi$.

Proof.

$$\begin{aligned}a + bi &= (a, 0) + (b, 0)(0, 1) \\ &= (a, 0) + (0, b) \\ &= (a, b).\end{aligned}$$

□

For $a, b \in \mathbb{R}$, we write $z = a + bi$; we call a and b the *real part* and *imaginary part* of z respectively, denoted by $a = \operatorname{Re}(z)$, $b = \operatorname{Im}(z)$; $\bar{z} = a - bi$ is called the *conjugate* of z .

Lemma 10.45 (Properties of conjugate). For $z, w \in \mathbb{C}$,

$$(i) \quad \overline{z + w} = \bar{z} + \bar{w}$$

$$(ii) \quad \overline{z\bar{w}} = \bar{z} w$$

$$(iii) \quad z + \bar{z} = 2 \operatorname{Re}(z), \quad z - \bar{z} = 2i \operatorname{Im}(z)$$

$$(iv) \quad z\bar{z} \in \mathbb{R} \text{ and } z\bar{z} \geq 0$$

For $z \in \mathbb{C}$, the *absolute value* of z is defined as

$$|z| := (z\bar{z})^{\frac{1}{2}}.$$

Lemma 10.46 (Properties of absolute value). For $z, w \in \mathbb{C}$,

$$(i) \quad |z| \geq 0$$

$$(ii) \quad |\bar{z}| = |z|$$

$$(iii) \quad |zw| = |z||w|$$

$$(iv) \quad |\operatorname{Re}(z)| \leq |z|$$

Proof.

(i) The square root is non-negative, by definition.

(ii) The conjugate of \bar{z} is z , and the rest follows by the definition of absolute value.

(iii) Let $z = a + bi$, $w = c + di$ where $a, b, c, d \in \mathbb{R}$. Then

$$\begin{aligned} |zw|^2 &= (ac - bd)^2 + (ad - bc)^2 \\ &= (a^2 + b^2)(c^2 + d^2) \\ &= |z|^2 |w|^2 = (|z||w|)^2 \end{aligned}$$

and the desired result follows by taking square roots on both sides.

(iv) Let $z = a + bi$. Note that $a^2 \leq a^2 + b^2$, hence

$$|\operatorname{Re}(z)| = |a| = \sqrt{a^2} \leq \sqrt{a^2 + b^2} = |z|.$$

□

Proposition 10.47 (Triangle inequality). For $z, w \in \mathbb{C}$,

$$|z + w| \leq |z| + |w|. \quad (10.1)$$

Proof. Let $z, w \in \mathbb{C}$. Note that the conjugate of $z\bar{w}$ is $\bar{z}w$, so $z\bar{w} + \bar{z}w = 2 \operatorname{Re}(z\bar{w})$. Hence

$$\begin{aligned} |z + w|^2 &= (z + w)(\overline{z + w}) = (z + w)(\bar{z} + \bar{w}) \\ &= z\bar{z} + z\bar{w} + \bar{z}w + w\bar{w} \\ &= |z|^2 + 2 \operatorname{Re}(z\bar{w}) + |w|^2 \\ &\leq |z|^2 + 2|z\bar{w}| + |w|^2 \\ &= |z|^2 + 2|z||w| + |w|^2 \\ &= (|z| + |w|)^2 \end{aligned}$$

and taking square roots yields the desired result. \square

Corollary 10.48 (Generalised triangle inequality). For $z_1, \dots, z_n \in \mathbb{C}$,

$$|z_1 + \dots + z_n| \leq |z_1| + \dots + |z_n|.$$

Proof. We have proven the case $n = 2$. Assume the statement holds for $n - 1$. Then

$$|z_1 + \dots + z_{n-1} + z_n| \leq |z_1 + \dots + z_{n-1}| + |z_n| \leq |z_1| + \dots + |z_n|,$$

which establishes the claim by induction. \square

Corollary 10.49. For $x, y, z \in \mathbb{C}$,

$$(i) \quad ||x| - |y|| \leq |x - y|;$$

$$(ii) \quad |x - y| \leq |x - z| + |z - y|.$$

Proof.

(i) By the triangle inequality,

$$|x| = |(x - y) + y| \leq |x - y| + |y|$$

so that

$$|x| - |y| \leq |x - y|.$$

Interchanging the roles of x and y in the above gives

$$|y| - |x| \leq |x - y|.$$

Hence

$$||x| - |y|| \leq |x - y|.$$

(ii) In the triangle inequality, replace x by $x - y$ and y by $y - z$.

\square

Proposition 10.50 (Cauchy–Schwarz inequality). *If $a_1, \dots, a_n, b_1, \dots, b_n \in \mathbb{C}$, then*

$$\left| \sum_{i=1}^n a_i \bar{b}_i \right|^2 \leq \sum_{i=1}^n |a_i|^2 \sum_{i=1}^n |b_i|^2. \quad (10.2)$$

Proof. For simplicity, we shall drop the upper and lower limits of the sums. Let

$$A = \sum |a_i|^2, \quad B = \sum |b_i|^2, \quad C = \sum a_i \bar{b}_i.$$

Then 10.2 becomes

$$|C|^2 \leq AB.$$

If $B = 0$, then $b_1 = \dots = b_n = 0$, and the conclusion is trivial. Now assume that $B > 0$. Then consider the sum

$$\begin{aligned} \sum |Ba_i - Cb_i|^2 &= \sum (Ba_i - Cb_i)(\overline{Ba_i - Cb_i}) \\ &= \sum (Ba_i - Cb_i)(B\bar{a}_i - \bar{C}\bar{b}_i) \\ &= B^2 \sum |a_i|^2 - B\bar{C} \sum a_i \bar{b}_i - BC \sum \bar{a}_i b_i + |C|^2 \sum |b_i|^2 \\ &= B^2 A - B|C|^2 \\ &= B(AB - |C|^2). \end{aligned}$$

Each term in $\sum |Ba_i - Cb_i|^2$ is non-negative, so $\sum |Ba_i - Cb_i|^2 \geq 0$. Thus

$$B(AB - |C|^2) \geq 0.$$

Since $B > 0$, it follows that $AB - |C|^2 \geq 0$, or $|C|^2 \leq AB$. This is the desired inequality.

(when does equality hold?) □

Define

$$\mathbb{C}^n = \{(z_1, \dots, z_n) \mid z_i \in \mathbb{C}\}.$$

We can define an inner product on \mathbb{C}^n : for $\mathbf{a}, \mathbf{b} \in \mathbb{C}^n$,

$$\langle \mathbf{a}, \mathbf{b} \rangle = \sum_{i=1}^n a_i \bar{b}_i.$$

We can also define the norm of $\mathbf{a} \in \mathbb{C}^n$:

$$|\mathbf{a}| = \langle \mathbf{a}, \mathbf{a} \rangle^{\frac{1}{2}}.$$

§10.4 Euclidean Space

For $n \in \mathbb{N}$, define

$$\mathbb{R}^n := \{(x_1, \dots, x_n) \mid x_i \in \mathbb{R}\}$$

where $\mathbf{x} = (x_1, \dots, x_n)$, x_i 's are called the coordinates of \mathbf{x} . The elements of \mathbb{R}^n are called *points*, or *vectors*.

Lemma 10.51. Let $\mathbf{x} = (x_1, \dots, x_n)$, $\mathbf{y} = (y_1, \dots, y_n)$. \mathbb{R}^n , with addition and scalar multiplication defined as

$$\begin{aligned}\mathbf{x} + \mathbf{y} &= (x_1 + y_1, \dots, x_n + y_n), \\ \alpha \mathbf{x} &= (\alpha x_1, \dots, \alpha x_n).\end{aligned}$$

for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $\alpha \in \mathbb{R}$, is a vector space over \mathbb{R} . Note that the zero element of \mathbb{R}^n is $\mathbf{0} = (0, \dots, 0)$.

Proof. These two operations satisfy the commutative, associative, and distributive laws (the proof is trivial, in view of the analogous laws for the real numbers). \square

We define the *inner product* of \mathbf{x} and \mathbf{y} by

$$\mathbf{x} \cdot \mathbf{y} := \sum_{i=1}^n x_i y_i,$$

and the *norm* of \mathbf{x} by

$$\|\mathbf{x}\| := \sqrt{\mathbf{x} \cdot \mathbf{x}}.$$

The structure now defined (the vector space \mathbb{R}^n with the above inner product and norm) is called the *Euclidean n -space*.

Lemma 10.52 (Basic properties of norm). Suppose $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^n$, $\alpha \in \mathbb{R}$.

(i) $\|\mathbf{x}\| \geq 0$, where equality holds if and only if $\mathbf{x} = \mathbf{0}$ (positive definiteness)

(ii) $\|\alpha \mathbf{x}\| = |\alpha| \|\mathbf{x}\|$ (homogeneity)

(iii) $\|\mathbf{x} \cdot \mathbf{y}\| \leq \|\mathbf{x}\| \|\mathbf{y}\|$ (Cauchy–Schwarz inequality)

(iv) $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ (triangle inequality)

(v) $\|\mathbf{x} - \mathbf{z}\| \leq \|\mathbf{x} - \mathbf{y}\| + \|\mathbf{y} - \mathbf{z}\|$ (triangle inequality)

Proof.

(i) Obvious from definition.

(ii) Obvious from definition.

(iii) We want to show

$$\sqrt{\sum_{i=1}^n x_i y_i} \leq \sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2},$$

or, squaring both sides,

$$\sum_{i=1}^n x_i y_i \leq \left(\sum_{i=1}^n x_i^2 \right) \left(\sum_{i=1}^n y_i^2 \right).$$

But this is simply the Cauchy–Schwarz inequality 10.2.

(iv) By (iii) we have

$$\begin{aligned} \|\mathbf{x} + \mathbf{y}\| &= (\mathbf{x} + \mathbf{y}) \cdot (\mathbf{x} + \mathbf{y}) \\ &= \mathbf{x} \cdot \mathbf{x} + 2\mathbf{x} \cdot \mathbf{y} + \mathbf{y} \cdot \mathbf{y} \\ &\leq \|\mathbf{x}\|^2 + 2\|\mathbf{x}\|\|\mathbf{y}\| + \|\mathbf{y}\|^2 \\ &= (\|\mathbf{x}\| + \|\mathbf{y}\|)^2. \end{aligned}$$

(v) This follows directly from (iv) by replacing \mathbf{x} by $\mathbf{x} - \mathbf{y}$, and \mathbf{y} by $\mathbf{y} - \mathbf{z}$.

□

Exercises

Exercise 10.1 ([Rud76] 1.1). If $r \in \mathbb{Q} \setminus \{0\}$ and $x \in \mathbb{R} \setminus \mathbb{Q}$, prove that $r + x \in \mathbb{R} \setminus \mathbb{Q}$ and $rx \in \mathbb{R} \setminus \mathbb{Q}$.

Solution. Prove by contradiction. If r and $r + x$ were both rational, then $x = (r + x) - r$ would also be rational. Similarly if rx were rational, then $x = \frac{rx}{r}$ would also be rational. \square

Exercise 10.2 ([Rud76] 1.2). Prove that there is no rational number whose square is 12.

Solution. Prove by contradiction. \square

Exercise 10.3 ([Rud76] 1.4). Let E be a nonempty subset of an ordered set; suppose α is a lower bound of E , and β is an upper bound of E . Prove that $\alpha \leq \beta$.

Solution. Since E is non-empty, there exists $x \in E$. By definition of lower and upper bounds, we have $\alpha \leq x \leq \beta$. \square

Exercise 10.4 ([Rud76] 1.8). Prove that no order can be defined in \mathbb{C} that turns it into an ordered field. *Hint:* -1 is a square.

Solution. By 10.14, an order $<$ that makes \mathbb{C} an ordered field would have to satisfy $-1 = i^2 > 0$, contradicting $1 > 0$. \square

Exercise 10.5 ([Rud76] 1.9, lexicographic order). Suppose $z = a + bi$, $w = c + di$. Define an order on \mathbb{C} as follows:

$$z < w \iff \begin{cases} a < c, \text{ or} \\ a = c, b < d. \end{cases}$$

Prove that this turns \mathbb{C} into an ordered set. Does this ordered set have the least upper bound property?

Solution. We show that this order turns \mathbb{C} into an ordered set.

- (i) Since the *real* numbers are ordered, we have $a < c$ or $a = c$ or $c < a$. In the first case $z < w$; in the third case $w < z$.

Now consider the second case where $a = c$. We must have $b < d$ or $b = d$ or $d < b$, which correspond to $z < w$, $z = w$, $w < z$ respectively.

Hence we have shown that either $z < w$ or $z = w$ or $w < z$.

- (ii) We now show that if $z < w$ and $w < u$, then $z < u$. Let $u = e + fi$.

Since $z < w$, we have either $a < c$, or $a = c$ and $b < d$. Since $w < u$, we have either $c < e$, or $c = e$ and $d < f$. Hence there are four possible cases:

- $a < c$ and $c < e$. Then $a < e$ and so $z < u$, as required.
- $a < c$ and $c = e$, and $d < f$. Again $a < e$, so $z < u$.
- $a = c$, and $b < d$ and $c < e$. Once again $a < e$ so $z < u$.
- $a = c$ and $b < d$, and $c = e$ and $d < f$. Then $a = e$ and $b < f$, so $z < u$.

\square

Exercise 10.6 ([Rud76] 1.10). Suppose $z = a + bi$, $w = u + iv$, and

$$a = \left(\frac{|w| + u}{2} \right)^{\frac{1}{2}}, \quad b = \left(\frac{|w| - u}{2} \right)^{\frac{1}{2}}.$$

Prove that $z^2 = w$ if $v \geq 0$ and that $\bar{z}^2 = w$ if $v \leq 0$. Conclude that every complex number (with one exception!) has two complex square roots.

Solution. We have

$$a^2 - b^2 = \frac{|w| + u}{2} - \frac{|w| - u}{2} = u,$$

and

$$2ab = (|w| + u)^{\frac{1}{2}} (|w| - u)^{\frac{1}{2}} = \left(|w|^2 - u^2 \right)^{\frac{1}{2}} = (v^2)^{\frac{1}{2}} = |v|.$$

Hence if $v \geq 0$,

$$z^2 = (a^2 - b^2) + 2abi = u + |v|i = w;$$

if $v \leq 0$,

$$\bar{z}^2 = (a^2 - b^2) - 2abi = u - |v|i = w.$$

Hence every non-zero w has two square roots $\pm z$ or $\pm \bar{z}$. Of course, 0 has only one square root, itself. \square

Exercise 10.7 ([Rud76] 1.11). If $z \in \mathbb{C}$, prove that there exists $r \geq 0$ and $w \in \mathbb{C}$ with $|w| = 1$ such that $z = rw$. Are w and r always uniquely determined by z ?

Solution. If $z = 0$, take $r = 0$ and $w = 1$; in this case w is not unique.

Otherwise take $r = |z|$ and $w = \frac{z}{|z|}$; these choices are unique, since if $z = rw$, we must have $r = r|w| = |rw| = |z|$ so $w = \frac{z}{r} = \frac{z}{|z|}$ are unique. \square

11 Basic Topology

Summary

- Metric space, subspace. Open ball, closed ball, boundedness. Open set, closed set. Interior, closure, boundary. Limit point.
- Compactness. Cantor intersection theorem, Heine–Borel theorem, Bolzano–Weierstrass theorem. Sequential compactness.
- Perfect sets. Cantor set.
- Connectedness.

Term	Notation
metric space	X, Y
metric	$d(p, q)$
general set	E
point in a set	p, q, r
open ball	$B_r(p)$
closed ball	$\overline{B}_r(p)$
punctured ball	$B'_r(p)$
neighbourhood	N
interior	E°
closure	\overline{E}
boundary	∂E
induced set	E'
compact set	K
open cover	\mathcal{U}
n -cell	I
Cantor set	C

Table 11.1: Notation for topological structures in Chapter 11

§11.1 Metric Spaces

Definitions and Examples

Definition 11.1 (Normed space). Let X be a vector space. A *norm* is a function $\|\cdot\| : X \rightarrow [0, \infty)$ if, for all $x, y \in X$ and constants α ,

(i) $\|x\| \geq 0$, where equality holds if and only if $x = 0$; (positive definiteness)

(ii) $\|\alpha x\| = |\alpha|\|x\|$; (homogeneity)

(iii) $\|x + y\| \leq \|x\| + \|y\|$. (triangle inequality)

A *normed space* $(V, \|\cdot\|)$ is a vector space V together with a norm $\|\cdot\|$.

Definition 11.2 (Metric space). Let X be a set. A *metric* is a function $d : X \times X \rightarrow [0, \infty)$ if, for all $p, q, r \in X$,

(i) $d(p, q) \geq 0$, where equality holds if and only if $p = q$; (positive definiteness)

(ii) $d(p, q) = d(q, p)$; (symmetry)

(iii) $d(p, q) \leq d(p, r) + d(r, q)$. (triangle inequality)

A *metric space* (X, d) is a set X together with a metric d .

For the rest of the chapter, X is taken to be a metric space, unless specified otherwise.

Lemma 11.3 (Norm induces metric). *Let X be a normed space. Then X is a metric space, with the metric $d(x, y) = \|x - y\|$ for every $x, y \in X$.*

Proof. Trivial; check the conditions for a metric. □

Example 11.4 (Metrics on \mathbb{R}^n). Each of the following functions define metrics on \mathbb{R}^n .

$$d_1(x, y) = \sum_{i=1}^n |x_i - y_i|;$$

$$d_2(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

$$d_\infty(x, y) = \max_{i \in \{1, 2, \dots, n\}} |x_i - y_i|.$$

These are called the ℓ^1 -, ℓ^2 - (or Euclidean) and ℓ^∞ -distances respectively.

The proof that each of d_1, d_2, d_∞ is a metric is mostly very routine, with the exception of proving that d_2 , the Euclidean distance, satisfies the triangle inequality. To establish this, recall that the Euclidean norm $\|x\|_2$ of a

vector $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ is

$$\|x\|_2 := \left(\sum_{i=1}^n x_i^2 \right)^{\frac{1}{2}} = \langle x, x \rangle^{\frac{1}{2}},$$

where the inner product is given by

$$\langle x, y \rangle := \sum_{i=1}^n x_i y_i.$$

Then $d_2(x, y) = \|x - y\|_2$, and so the triangle inequality is the statement that

$$\|w - y\|_2 \leq \|w - x\|_2 + \|x - y\|_2.$$

This follows immediately by taking $u = w - x$ and $v = x - y$ in the following lemma.

Lemma. *If $u, v \in \mathbb{R}^n$ then $\|u + v\|_2 \leq \|u\|_2 + \|v\|_2$.*

Proof. Since $\|u\|_2 \geq 0$ for all $u \in \mathbb{R}^n$, squaring both sides of the desired inequality gives

$$\|u + v\|_2^2 \leq \|u\|_2^2 + 2\|u\|_2\|v\|_2 + \|v\|_2^2.$$

But since

$$\|u + v\|_2^2 = \langle u + v, u + v \rangle = \|u\|_2^2 + 2\langle u, v \rangle + \|v\|_2^2,$$

this inequality is immediate from the Cauchy–Schwarz inequality, that is to say the inequality

$$|\langle u, v \rangle| \leq \|u\|_2\|v\|_2.$$

□

A metric space (X, d) naturally induces a metric on any of its subsets.

Definition 11.5 (Subspace). Suppose (X, d) is a metric space, $Y \subset X$. Then the restriction of d to $Y \times Y$ gives Y a metric so that $(Y, d_{Y \times Y})$ is a metric space. We call Y equipped with this metric a *subspace*.

Balls and Boundedness

Definition 11.6 (Balls).

- (i) The **open ball** centred at $p \in X$ with radius $r > 0$ is the set

$$B_r(p) := \{q \in X \mid d(p, q) < r\}.$$

- (ii) The **closed ball** centred at p with radius r is

$$\overline{B}_r(p) := \{q \in X \mid d(p, q) \leq r\}.$$

- (iii) The **punctured ball** is the open ball excluding its centre:

$$B'_r(p) := \{q \in X \mid 0 < d(p, q) < r\}.$$

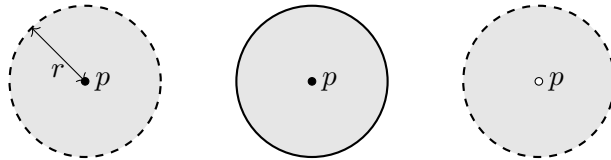


Figure 11.1: Open ball, closed ball, punctured ball

Example 11.7. Considering \mathbb{R}^3 with the Euclidean metric, $B_1(0)$ really is what we understand geometrically as a ball (minus its boundary, the unit sphere), whilst $\overline{B}_1(0)$ contains the unit sphere and everything inside it.

Remark. We caution that this intuitive picture of the closed ball being the open ball “together with its boundary” is totally misleading in general. For instance, in the discrete metric on a set X , the open ball $B_1(a)$ contains only the point a , whereas the closed ball $\overline{B}_1(a)$ is the whole of X .

Definition 11.8 (Bounded). $E \subset X$ is said to be **bounded** if E is contained in some open ball; that is, there exists $M \in \mathbb{R}$ and $p \in X$ such that $E \subset B_M(p)$.

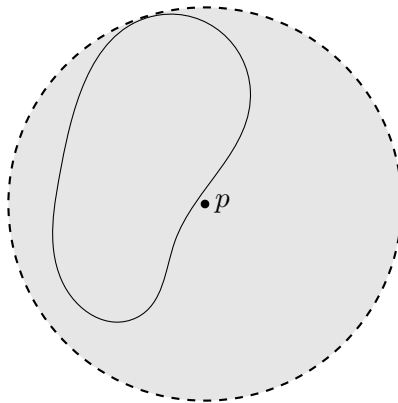


Figure 11.2: Bounded set

Proposition 11.9. *Let $E \subset X$. Then the following are equivalent:*

- (i) E is bounded;
- (ii) E is contained in some closed ball;
- (iii) The set $\{d(x, y) \mid x, y \in E\}$ is a bounded subset of \mathbb{R} .

Proof.

$(i) \implies (ii)$ This is obvious.

$(ii) \implies (iii)$ This follows immediately from the triangle inequality.

$(iii) \implies (i)$ Suppose E satisfies (iii), then there exists $r \in \mathbb{R}$ such that $d(x, y) \leq r$ for all $x, y \in E$. If $E = \emptyset$, then E is certainly bounded. Otherwise, let $p \in E$ be an arbitrary point. Then $E \subset B_{r+1}(p)$. \square

Open and Closed Sets

We say $N \subset X$ is a **neighbourhood** of $p \in X$ if there exists $\varepsilon > 0$ such that $B_\varepsilon(p) \subset N$.

Definition 11.10 (Open set). $E \subset X$ is **open** (in X) if it is a neighbourhood of all its elements; that is, for all $p \in E$, there exists $\varepsilon > 0$ such that $B_\varepsilon(p) \subset E$.

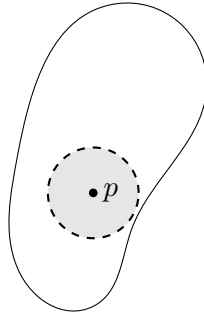


Figure 11.3: Open set

Lemma 11.11. Any open ball is open.

Proof. Let $B_r(p)$ be an open ball.

Let $q \in B_r(p)$. To show that $B_r(p)$ is open, we will show that $B_\varepsilon(q) \subset B_r(p)$ for some $\varepsilon > 0$.

Take $\varepsilon = r - d(p, q)$. [Note that $\varepsilon > 0$.] Let $s \in B_\varepsilon(q)$. By the triangle inequality,

$$\begin{aligned} d(p, s) &\leq d(q, s) + d(p, q) \\ &< \varepsilon + d(p, q) = r \end{aligned}$$

so $s \in B_r(p)$, which implies $B_\varepsilon(q) \subset B_r(p)$. □

Lemma 11.12.

- (i) Both \emptyset and X are open.
- (ii) For any indexing set I and collection of open sets $\{E_i \mid i \in I\}$, $\bigcup_{i \in I} E_i$ is open.
- (iii) For any finite indexing set I and collection of open sets $\{E_i \mid i \in I\}$, $\bigcap_{i \in I} E_i$ is open.

Proof.

- (i) Obvious by definition.
- (ii) If $p \in \bigcup_{i \in I} E_i$, then $p \in E_i$ for some $i \in I$. Since E_i is open, there exists $\varepsilon > 0$ such that $B_\varepsilon(p) \subset E_i$ and hence $B_\varepsilon(p) \subset \bigcup_{i \in I} E_i$.
- (iii) Suppose that I is finite and that $p \in \bigcap_{i \in I} E_i$. For each $i \in I$, we have $p \in E_i$ and so there exists δ_i such that $B_{\delta_i}(p) \subset E_i$. Set $\delta = \min_{i \in I} \delta_i$, then $\delta > 0$ (here it is, of course, crucial that I be finite), and $B_\delta(p) \subset B_{\delta_i}(p) \subset E_i$ for all i . Therefore $B_\delta(p) \subset \bigcap_{i \in I} E_i$.

□

Remark. While the indexing set I in (ii) can be arbitrary, the indexing set in (iii) must be finite. For instance, $E_n = \left(-\frac{1}{n}, \frac{1}{n}\right)$ are open in \mathbb{R} , but their intersection $\bigcap_{n=1}^{\infty} E_n = \{0\}$ is not open.

Suppose Y is a subspace of X . We say that E is *open relative to Y* if for all $p \in E$, there exists $\varepsilon > 0$ such that $B_\varepsilon(p) \cap Y \subset E$. (Note that $B_\varepsilon(p) \cap Y$ is in the open ball in Y^1 , because the metric $d' : Y \times Y \rightarrow \mathbb{R}$ is the restriction to $Y \times Y$ of the metric $d : X \times X \rightarrow \mathbb{R}$ on X .)

Proposition 11.13. *Suppose Y is a subspace of X , $E \subset Y$. Then E is open relative to Y if and only if there exists an open subset G of X such that $E = Y \cap G$.*

Proof.

⟹ We prove by construction; that is, construct the required set G .

Suppose E is open relative to Y . For each $p \in E$, by openness of E , there exists $r_p > 0$ such that $B_{r_p}(p) \cap Y \subset E$. Consider the union

$$\bigcup_{p \in E} (B_{r_p}(p) \cap Y) \subset E.$$

Note that we can write

$$\bigcup_{p \in E} (B_{r_p}(p) \cap Y) = \left(\bigcup_{p \in E} B_{r_p}(p) \right) \cap Y \subset E.$$

Let

$$G = \bigcup_{p \in E} B_{r_p}(p),$$

then we have $G \cap Y \subset E$.

Since G is an intersection of open balls (which are open sets), by 11.12, G is an open subset of X .

Note for each $p \in E \subset Y$, we have $p \in Y$, and $p \in B_{r_p}(p)$ for some $r_p > 0$, so $p \in \bigcup_{p \in E} B_{r_p}(p) = G$. Hence $p \in G \cap Y$. This shows $E \subset G \cap Y$.

Hence $E = G \cap Y$.

⟸ Suppose $E = G \cap Y$ for some open subset G of X .

Let $p \in E$. Since $p \in G$, by the openness of G , there exists $r_p > 0$ such that $B_{r_p}(p) \subset G$. Then $B_{r_p}(p) \cap Y \subset G \cap Y = E$. Thus by definition E is open relative to Y . □

The complement of an open set is a closed set.

Definition 11.14 (Closed set). $E \subset X$ is **closed** if its complement $E^c = X \setminus E$ is open.

Lemma 11.15. *Any closed ball is closed.*

Proof. To prove that $\overline{B}_r(p)$ is closed, we need to show that its complement

$$\overline{B}_r(p)^c = \{q \in X \mid d(p, q) > r\}$$

¹notice that the definition of an open ball depends on the metric space!

is open.

Let $s \in \overline{B_r(p)}^c$. Take $\varepsilon > 0$ such that $r + \varepsilon < d(p, s)$; that is, $\varepsilon < d(p, s) - r$.

Let $q \in B_\varepsilon(s)$, then $d(q, s) < \varepsilon$. Thus $d(q, s) < d(p, s) - r$, or $r < d(p, s) - d(q, s)$. Then by the triangle inequality,

$$\begin{aligned} d(p, q) &\geq d(p, s) - d(q, s) \\ &> r \end{aligned}$$

Hence $q \in \overline{B_r(p)}^c$, and so $B_\varepsilon(s) \subset \overline{B_r(p)}^c$. Therefore $\overline{B_r(p)}^c$ is open, so $\overline{B_r(p)}$ is closed. \square

Lemma 11.16.

- (i) Both \emptyset and X are closed.
- (ii) For any indexing set I and collection of closed sets $\{F_i \mid i \in I\}$, $\bigcap_{i \in I} F_i$ is closed.
- (iii) For any finite indexing set I and collection of closed sets $\{F_i \mid i \in I\}$, $\bigcup_{i \in I} F_i$ is closed.

Proof. From 11.12, simply take complements and apply de Morgan's laws. \square

Remark. The indexing set in (iii) must be finite; for instance, the closed intervals $F_n = \left[-1 + \frac{1}{n}, 1 - \frac{1}{n}\right]$ are all closed in \mathbb{R} , but their union $\bigcup_{n=1}^{\infty} F_n = (-1, 1)$ is open.

Interior, Closure, Boundary

Definition 11.17. Suppose $E \subset X$.

- (i) The **interior** E° of the set E is the union of all open subsets of X contained in E ; we call $p \in E^\circ$ an **interior point** of E .
- (ii) The **closure** \bar{E} of the set E is the intersection of all closed subsets of X containing E .
We say E is **dense** if $\bar{E} = X$.
- (iii) The **boundary** of E is $\partial E = \bar{E} \setminus E^\circ$; we call $p \in \partial E$ a **boundary point** of E .

In the figure below, the black outline represents the boundary; the grey area within represents the interior; the union represents the closure.

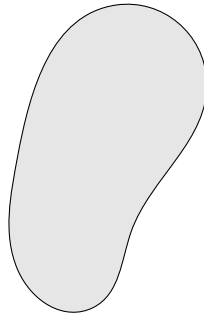


Figure 11.4: Interior, closure, boundary

Example 11.18.

- The interior of the closed interval $[a, b]$ is the open interval (a, b) .
- \mathbb{Q} is dense in \mathbb{R} .

Remark. E and E° do not necessarily have the same closures; for example, take $E = \mathbb{Q}$, then $\bar{E} = \mathbb{R}$ and $\overline{E^\circ} = \emptyset$.

Likewise, E and \bar{E} do not necessarily have the same interiors; for example, take $E = (-1, 0) \cup (0, 1) \subset \mathbb{R}$. Then $E^\circ = (-1, 0) \cup (0, 1)$ and $(\bar{E})^\circ = [-1, 1]$.

Lemma 11.19. Suppose $E \subset X$.

- (i) E is open if and only if $E = E^\circ$.
(That is, E is open if and only if every point of E is an interior point.)
- (ii) E is closed if and only if $E = \bar{E}$.

Proof.

- (i) $\boxed{\implies}$ Suppose E is open. By assumption, E is an open subset of X contained in E (since $E \subset E$), so $E \subset E^\circ$.

We now show the opposite containment. Let $x \in E^\circ$. Then x is in some open subset of X contained in E , so $x \in E$. Hence $E^\circ \subset E$.

Therefore $E = E^\circ$.

$\boxed{\Leftarrow}$ Since an arbitrary union of open sets is open, E° is open. Since $E = E^\circ$, we have that E is open.

(ii) $\boxed{\Rightarrow}$ Suppose E is closed. Then $E \subset \overline{E}$.

We now show the opposite containment. Let $x \in \overline{E}$. Then x is in every closed subset of X containing E , so $x \in E$. Hence $x \in E$.

Therefore $E = \overline{E}$.

$\boxed{\Leftarrow}$ Since an arbitrary intersection of closed sets is closed, \overline{E} is closed. Since $E = \overline{E}$, we have that E is closed.

□

Proposition 11.20. *Suppose $E \subset X$. Then $p \in \overline{E}$ if and only if every open ball centred at p contains a point of E .*

Proof.

$\boxed{\Rightarrow}$ Let $p \in \overline{E}$.

Suppose, for a contradiction, that there exists an open ball $B_\varepsilon(p)$ that does not meet E . Then $B_\varepsilon(p)^c$ is a closed set containing E . Therefore $B_\varepsilon(p)^c$ contains \overline{E} , and hence it contains p , which is obviously nonsense.

$\boxed{\Leftarrow}$ Suppose that every ball $B_\varepsilon(p)$ meets E .

Suppose, for a contradiction, that $p \notin \overline{E}$. Since \overline{E}^c is open, there is a ball $B_\varepsilon(p)$ contained in \overline{E}^c , and hence in E^c , contrary to assumption. □

Remark. A particular consequence of this is that $E \subset X$ is dense if and only if it meets every open set in X .

Lemma 11.21 (Properties of closure and interior). *Suppose $A, B \subset X$. Then*

(i) $\overline{A \cup B} = \overline{A} \cup \overline{B}$

(ii) $\overline{A \cap B} \subset \overline{A} \cap \overline{B}$

(iii) $(A \cup B)^\circ \supset A^\circ \cup B^\circ$

(iv) $(A \cap B)^\circ = A^\circ \cap B^\circ$

(v) $(A^\circ)^c = \overline{A^c}$

(vi) $(\overline{A})^c = (A^c)^\circ$

Limit Points

Definition 11.22.

- (i) $p \in X$ (not necessarily in E) is an **adherent point** of E (or is *adherent* to E) if $B_\varepsilon(p) \cap E \neq \emptyset$ for all $\varepsilon > 0$.
- (ii) $p \in X$ is a **limit point** of E if, for all $\varepsilon > 0$, there exists $q \in E \setminus \{p\}$ such that $q \in B_\varepsilon(p)$. (In other words, p is a limit point of E if and only if p adheres to $E \setminus \{p\}$.)
- The **induced set** of E , denoted by E' , is the set of all limit points of E in X .
- (iii) $p \in E$ is an **isolated point** of E if p is not a limit point of E (that is, there exists $\varepsilon > 0$ such that $B_\varepsilon(p) \cap E = \{p\}$).

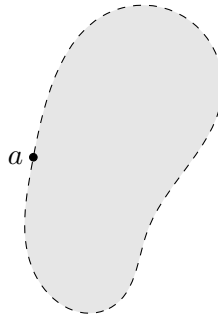


Figure 11.5: Adherent point, limit point, isolated point

Example 11.23 (Adherent point).

- If $p \in E$, then p adheres to E because every ball contains p .
- If $E \subset \mathbb{R}$ is bounded above, then $\sup E$ is adherent to E .

Example 11.24 (Limit point).

- The set $\left\{ \frac{1}{n} \mid n \in \mathbb{N} \right\}$ has 0 as a limit point.
- The set of rational numbers has every real number as a limit point.
- Every point of $[a, b]$ is a limit point of the set of numbers in (a, b) .
- Consider \mathbb{R}^2 . The set of limit points of any open ball $B_r(p)$ is the closed ball $\overline{B}_r(p)$, which is also the closure of $B_r(p)$.
- Consider $\mathbb{Q} \subset \mathbb{R}$. $\mathbb{Q}' = \overline{\mathbb{Q}} = \mathbb{R}$.

Proposition 11.25. *If p is a limit point of E , then every open ball of p contains infinitely many points of E .*

Proof. Suppose, for a contradiction, that there exists an open ball $B_r(p)$ which contains only a finite number of points of E distinct from p ; let

$$B_r(p) = \{q_1, \dots, q_n\},$$

where $p \neq q_i$ for $i = 1, \dots, n$. Take

$$r = \min\{d(p, q_1), \dots, d(p, q_n)\},$$

then $B_r(p)$ contains no points of E distinct from p , which is a contradiction. \square

Corollary 11.26. *A finite point set has no limit points.*

Remark. The converse is not true; for example, \mathbb{N} is an infinite set with no limit points. In a later section we will show that infinite sets contained in some open ball always have a limit point; this result is known as the Bolzano–Weierstrass theorem (11.49).

A closed set was defined to be the complement of an open set. The next result characterises closed sets in another way.

Lemma 11.27. *Suppose $E \subset X$. Then E is closed if and only if it contains all its limit points.*

Proof.

\Rightarrow Suppose E is closed. Let p be a limit point of E . We want to show $p \in E$.

Suppose, for a contradiction, that $p \notin E$. Then $p \in E^c$. Since E^c is open, there exists $\varepsilon > 0$ such that $B_\varepsilon(p) \subset E^c$. Thus $B_\varepsilon(p)$ contains no points of E , contradicting the fact that p is a limit point of E .

\Leftarrow Suppose E contains all its limit points. To show that E is closed, we want to show that E^c is open.

Let $p \in E^c$. Then p is not a limit point of E , so there exists some ball $B_\varepsilon(p)$ which does not intersect E , so $B_\varepsilon(p) \subset E^c$. Hence E^c is open, so E is closed. \square

Lemma 11.28. *Suppose $E \subset X$. Then E' is a closed subset of X .*

Proof. To prove that E' is closed, we will show its complement $(E')^c$ is open.

Let $p \in (E')^c$. Then $p \notin E'$, so p is not a limit point of E ; thus, there exists a ball $B_\varepsilon(p)$ whose intersection with E is either empty or $\{p\}$ (depending on whether $p \in E$ or not).

We will show that $B_{\frac{\varepsilon}{2}}(p) \subset (E')^c$. Let $q \in B_{\frac{\varepsilon}{2}}(p)$.

Case 1: $q = p$. Then clearly $q \in (E')^c$.

Case 2: $q \neq p$. There is some ball about q which is contained in $B_\varepsilon(p)$, but does not contain p : the ball $B_\delta(q)$ where $\delta = \min\left(\frac{\varepsilon}{2}, d(p, q)\right)$ has this property. This ball meets E in the empty set, and so $q \in (E')^c$ in this case too.

\square

The next result provides a useful expression for the closure of a set; it states that every point of \overline{E} is either a limit point of E , or in E .

Lemma 11.29. *Suppose $E \subset X$. Then $\overline{E} = E \cup E'$.*

Proof. We show double inclusion.

- $E \cup E' \subset \overline{E}$ Obviously $E \subset \overline{E}$, so we need only show that $E' \subset \overline{E}$.

We prove by contrapositive. Suppose $p \in \overline{E}^c$. Since \overline{E}^c is open, there is some ball $B_\varepsilon(p)$ which lies in \overline{E}^c , and hence also in E^c , and therefore p cannot be a limit point of E .

- $\overline{E} \subset E \cup E'$ If $p \in \overline{E}$, we saw in Lemma 5.1.5 that there is a sequence (x_n) of elements of E with $x_n \rightarrow p$. If $x_n = p$ for some n then we are done, since this implies that $p \in E$. Suppose, then, that $x_n \neq p$ for all n . Let $\varepsilon > 0$ be given, for sufficiently large n , all the x_n are elements of $B_\varepsilon(p) \setminus \{p\}$, and they all lie in E . It follows that p is a limit point of E , and so we are done in this case also. to do

□

Lemma 11.30. Suppose $E \subset X$. Then \overline{E} is the smallest closed set containing E .

Proof. Let $F \supset E$ be some closed set in X . We will show that $\overline{E} \subset F$.

Let p be a limit point of E . Then p is a limit point of F . But since F is closed, by 11.27, F contains all its limit points, so all the limit points of E are in F . Hence $\overline{E} \subset F$. to do

□

Lemma 11.31. Suppose non-empty $E \subset \mathbb{R}$ is bounded above. Let $y = \sup E$. Then $y \in \overline{E}$. Hence $y \in E$ if E is closed.

Proof. If $y \in E$, since $E \subset \overline{E}$ we have that $y \in \overline{E}$.

For the second part, assume $y \notin E$. For every $h > 0$ there exists then a point $x \in E$ such that $y - h < x < y$, for otherwise $y - h$ would be an upper bound of E . Thus y is a limit point of E . Hence $y \in \overline{E}$. review proof

□

§11.2 Compactness

Definitions and Properties

Definition 11.32 (Open cover). An **open cover** of $K \subset X$ is a collection of open sets $\mathcal{U} = \{U_i \mid i \in I\}$ such that

$$K \subset \bigcup_{i \in I} U_i.$$

A **subcover** of \mathcal{U} is a subcollection $\{U_i \mid i \in I'\}$, where $I' \subset I$, which is an open cover of K . If I' is finite, then it is called a **finite subcover**.

Definition 11.33 (Compactness). $K \subset X$ is **compact** if every open cover of K contains a finite subcover.

That is, if $\mathcal{U} = \{U_i \mid i \in I\}$ is an open cover of K , then there are finitely many indices $i_1, \dots, i_n \in I$ such that

$$K \subset \bigcup_{k=1}^n U_{i_k}.$$

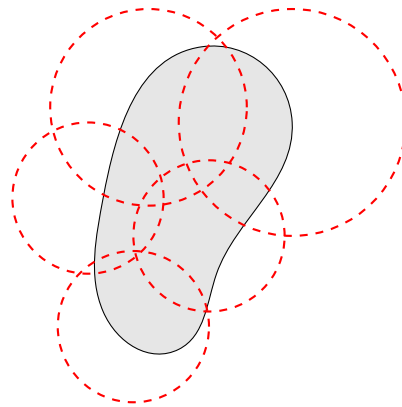


Figure 11.6: Compact set

Example 11.34.

- \mathbb{R} is not compact; for instance, the open cover $\{(-n, n) \mid n \in \mathbb{N}\}$ has no finite subcover.
- \mathbb{Z} is not compact in \mathbb{R} ; for instance, the open cover $\left\{ \left(n - \frac{1}{2}, n + \frac{1}{2} \right) \mid n \in \mathbb{Z} \right\}$ has no finite subcover.
- $[0, 1]$ is compact. (See 11.39 for the proof.)

Lemma 11.35. Every finite set is compact.

Proof. Let $E = \{p_1, \dots, p_n\}$. Let $\mathcal{U} = \{U_i \mid i \in I\}$ be an open cover of E . We will construct a finite subcover of E .

For each point $p_k \in E$, choose *one* U_{i_k} such that $p_k \in U_{i_k}$. Then $\{U_{i_k} \mid k = 1, \dots, n\}$ is a finite subcover of \mathcal{U} . □

Notice earlier than if $E \subset Y \subset X$, then E may be open relative to Y , but not open relative to X ; this implies that the property of being open depends on the space in which E is embedded. Compactness, however, behaves better, as shown in the next result; it is independent of the metric space.

Proposition 11.36. *Suppose Y is a subspace of X , and $K \subset Y$. Then K is compact relative to X if and only if K is compact relative to Y .*

Proof.

\implies Suppose K is compact relative to X .

Let \mathcal{U} be an open cover of K in Y ; that is, $\mathcal{U} = \{U_i \mid i \in I\}$ is a collection of sets open relative to Y , such that $K \subset \bigcup_{i \in I} U_i$. We want to show that \mathcal{U} has a finite subcover.

Since each U_i is open relative to Y , by 11.13, there exists V_i open relative to X such that $U_i = Y \cap V_i$. Consider the open cover $\{V_i \mid i \in I\}$ of K . Since K is compact relative to X , there exist finitely many indices i_1, \dots, i_n such that

$$K \subset \bigcup_{k=1}^n V_{i_k}.$$

Since $K \subset \bigcup_{k=1}^n V_{i_k}$ and $K \subset Y$, we have that

$$K \subset \left(\bigcup_{k=1}^n V_{i_k} \right) \cap Y = \bigcup_{k=1}^n (Y \cap V_{i_k}) = \bigcup_{k=1}^n U_{i_k},$$

where $\{U_{i_k} \mid k = 1, \dots, n\}$ forms a finite subcover of \mathcal{U} . Hence K is compact relative to Y .

\impliedby Suppose K is compact relative to Y . Let \mathcal{V} be an open cover of K in X ; that is, $\mathcal{V} = \{V_i \mid i \in I\}$ is a collection of open subsets of X which covers K . We want to show that \mathcal{V} has a finite subcover.

For $i \in I$, let $U_i = Y \cap V_i$. Then $\{U_i \mid i \in I\}$ cover K in Y . By compactness of K in Y , there exist finitely many indices i_1, \dots, i_n such that

$$K \subset \bigcup_{k=1}^n U_{i_k} \subset \bigcup_{k=1}^n V_{i_k}$$

since $U_i \subset V_i$. □

Proposition 11.37. *Compact subsets of metric spaces are bounded.*

Proof. Suppose $K \subset X$ is compact. To prove that K is bounded, we want to construct some open ball that contains the entirety of K .

Fix $p \in K$. For $n \in \mathbb{N}$, let $U_n = B_n(p)$. Then $\{U_n \mid n \in \mathbb{N}\}$ is an open cover of K . By compactness of K , there exists a finite subcover

$$\{U_{n_i} \mid i = 1, \dots, m\}.$$

But note that $U_{n_1} \subset \dots \subset U_{n_m}$, so U_{n_m} contains K . Hence K is bounded. □

Proposition 11.38. *Compact subsets of metric spaces are closed.*

Proof. Let $K \subset X$ be compact. To prove that K is closed, we need to show that K^c is open. Let $p \in K^c$; our goal is to show that there exists $\varepsilon > 0$ such that $B_\varepsilon(p) \subset K^c$, or $B_\varepsilon(p) \cap K = \emptyset$.

For all $q_i \in K$, consider the pair of open balls $B_{r_i}(p)$ and $B_{r_i}(q_i)$, where $r_i < \frac{1}{2}d(p, q_i)$. Since K is compact, there exists finite many points $q_{i_1}, \dots, q_{i_n} \in K$ such that

$$K \subset \bigcup_{k=1}^n B_{r_{i_k}}(q_{i_k}) = W.$$

Consider the intersection

$$\bigcap_{k=1}^n B_{r_{i_k}}(p),$$

which is an open ball at p of radius $\min\{d(p, q_{i_k}) \mid k = 1, \dots, n\}$.

Claim. $\varepsilon = \min\{d(p, q_{i_k}) \mid k = 1, \dots, n\}$.

Note that $B_\varepsilon(p) \subset B_{r_{i_k}}(p)$ for all $k = 1, \dots, n$. By construction, for all $q_i \in K$, the open balls $B_{r_i}(p)$ and $B_{r_i}(q_i)$ are disjoint. In particular,

$$B_\varepsilon(p) \cap B_{r_{i_k}}(q_{i_k}) = \emptyset \quad (k = 1, \dots, n)$$

Then

$$B_\varepsilon(p) \cap W = B_\varepsilon(p) \cap \left(\bigcup_{k=1}^n B_{r_{i_k}}(q_{i_k}) \right) = \bigcup_{k=1}^n \left(B_\varepsilon(p) \cap B_{r_{i_k}}(q_{i_k}) \right) = \emptyset$$

as desired. □

Proposition 11.39. *Closed subsets of compact sets are compact.*

Proof. Suppose $K \subset X$ is compact, $F \subset K$ is closed (relative to X). We will show that F is compact. Let $\mathcal{U} = \{U_i \mid i \in I\}$ be an open cover of F . We will construct a finite subcover of \mathcal{U} .

Since F is closed, its complement F^c is open. Consider the union

$$\Omega = \mathcal{U} \cup \{F^c\},$$

which is an open cover of K .

Since K is compact, there exists a finite subcover of Ω , given by

$$\Phi = \{U_{i_1}, \dots, U_{i_n}, F^c\}$$

which covers K , and hence F . Now remove F^c from Φ to obtain

$$\Phi' = \{U_{i_1}, \dots, U_{i_n}\},$$

which is an open cover of F , since $F^c \cap F = \emptyset$. Hence Φ' is a finite subcover of \mathcal{U} , so F is compact. □

Remark. Caution: this does *not* say “closed sets are compact”! In fact, closed sets are not necessarily compact. For instance, \mathbb{R} is closed in \mathbb{R} , but it is not compact because it is not bounded.

Note that closed and bounded sets are not necessarily compact for general metric spaces, but they are compact in \mathbb{R}^n (by 11.48).

Corollary 11.40. *If F is closed and K is compact, then $F \cap K$ is compact.*

Proof. Suppose F is closed, K is compact. By 11.38, K is closed. By 11.16, the intersection of two closed sets is closed, so $F \cap K$ is closed.

Since $F \cap K \subset K$ is a closed subset of a compact set K , by 11.39, $F \cap K$ is compact. □

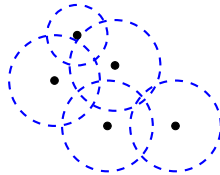
Heine–Borel Theorem

Proposition 11.41. *K is compact if and only if every infinite subset of K has a limit point in K .*

Proof.

\implies Suppose K is compact. Let E be an infinite subset of K . Suppose, for a contradiction, that E has no limit point in K .

For all $p \in K$, p is not a limit point of E , so there exists $r_p > 0$ such that $B_{r_p}(p) \cap E \setminus \{p\} = \emptyset$.



Consider the open cover of K given by the collection of open balls at each $p \in K$:

$$\mathcal{U} = \{B_{r_p}(p) \mid p \in E\}.$$

It is clear that \mathcal{U} has no finite subcover, since E is infinite, and each $B_{r_p}(p)$ contains at most one point of E .

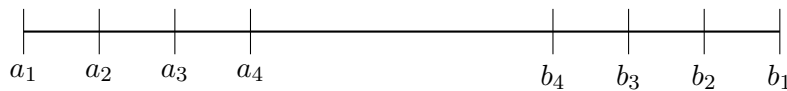
Since $E \subset K$, the above is also true for K . This contradicts the compactness of K .

\impliedby Suppose every infinite subset of K that has a limit point in K . Fix an arbitrary open cover $\mathcal{U} = \{U_i \mid i \in I\}$ of K . We will show that \mathcal{U} has a finite subcover, by construction.

Before that, we will reindex \mathcal{U} to make it more convenient, as follows. By the definition of a cover, every $p \in K$ is contained in some U_i . Pick *one* such U_i for each $p \in K$, and call it U_p . Then our open cover is now $\mathcal{U} = \{U_p \mid p \in K\}$, and for all $p \in K$ we have $p \in U_p$.

To complete proof

Proposition 11.42 (Nested interval theorem). *Suppose (I_n) is a decreasing sequence of closed and bounded intervals in \mathbb{R} ; that is, $I_1 \supset I_2 \supset \dots$. Then*

$$\bigcap_{n=1}^{\infty} I_n \neq \emptyset.$$


Proof. Let $I_n = [a_n, b_n]$, for $n = 1, 2, \dots$

Let $E = \{a_n \mid n \in \mathbb{N}\}$. Since E is non-empty and bounded above (by b_1), it has a supremum in \mathbb{R} ; let $x = \sup E$.

Claim. $x \in \bigcap_{n=1}^{\infty} I_n$.

Since x is the supremum, we have that $a_n \leq x$ for all $n \in \mathbb{N}$. Note that for $m > n$, $I_n \supset I_m$ implies $a_n \leq a_m \leq b_m \leq b_n$. This means b_n is an upper bound for all a_n ; hence $x \leq b_n$ for all $n \in \mathbb{N}$.

Therefore $x \in I_n$ for $n = 1, 2, \dots$ □

To generalise the notion of intervals, we define a k -cell as

$$\{(x_1, \dots, x_k) \in \mathbb{R}^k \mid a_i \leq x_i \leq b_i, 1 \leq i \leq k\}.$$

Example 11.43. A 1-cell is an interval, a 2-cell is a rectangle, and a 3-cell is a rectangular solid. In this regard, we can think of a k -cell as a higher-dimensional version of a rectangle or rectangular solid; it is the Cartesian product of k closed intervals.

The previous result can be generalised to k -cells, which we will now prove.

Proposition 11.44. *Suppose (I_n) is a decreasing sequence of k -cells; that is, $I_1 \supset I_2 \supset \dots$. Then $\bigcap_{n=1}^{\infty} I_n \neq \emptyset$.*

Proof. Let I_n consist of all points $\mathbf{x} = (x_1, \dots, x_k)$ such that

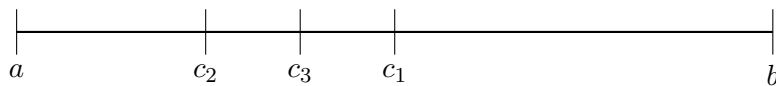
$$a_{n,i} \leq x_i \leq b_{n,i} \quad (1 \leq i \leq k; n = 1, 2, \dots),$$

and put $I_{n,i} = [a_{n,i}, b_{n,i}]$. For each i , the sequence $(I_{n,i})$ satisfies the hypotheses of 11.42. Hence there are real numbers x'_i ($1 \leq i \leq k$) such that

$$a_{n,i} \leq x'_i \leq b_{n,i} \quad (1 \leq i \leq k; n = 1, 2, \dots).$$

Setting $\mathbf{x}' = (x'_1, \dots, x'_k)$, we see that $\mathbf{x}' \in I_n$ for $n = 1, 2, \dots$. Hence $\bigcap_{n=1}^{\infty} I_n \neq \emptyset$, as desired. □

Lemma 11.45. *Every closed interval is compact (in \mathbb{R}).*



Proof. Suppose, for a contradiction, that a closed interval $[a, b] \subset \mathbb{R}$ is not compact. Then there exists an open cover $\mathcal{U} = \{U_i \mid i \in I\}$ with no finite subcover.

Let $c_1 = \frac{1}{2}(a, b)$. Subdivide $[a, b]$ into subintervals $[a, c_1]$ and $[c_1, b]$. Then \mathcal{U} covers $[a, c_1]$ and $[c_1, b]$, but at least one of these subintervals has no finite subcover (if not, then both subintervals have finite subcovers, so we can take the union of the two finite subcovers to obtain a larger subcover of the entire interval). WLOG, assume $[a, c_1]$ has no finite subcover; let $I_1 = [a, c_1]$.

Again subdivide I_1 in half to get $[a, c_2]$ and $[c_2, c_1]$. At least one of these subintervals has no finite subcover.

Repeat the above process of subdividing intervals into half. Then we obtain a decreasing sequence of closed intervals

$$I_1 \supset I_2 \supset I_3 \supset \dots$$

where all of them have no finite subcover of \mathcal{U} .

By the nested interval theorem (11.42), there exists $x' \in I_n$ for all $n \in \mathbb{N}$. Notice x' is in some U_i , which is open. Then there exists $\varepsilon > 0$ such that $B_\varepsilon(x') \subset U_i$.

Since the length of the subintervals is decreasing and tends to zero, there exists some subinterval I_n so small such that $I_n \subset B_\varepsilon(x')$. This means $I_n \subset U_i$, so U_i itself is an open cover of I_n , which contradicts the fact that I_n has no finite subcover of \mathcal{U} . \square

We now show a more general result.

Lemma 11.46. *Every k -cell is compact (in \mathbb{R}^k).*

Proof. We proceed in a similar manner to the proof the previous result.

Suppose I is a k -cell; that is,

$$I = \{(x_1, \dots, x_k) \mid a_i \leq x_i \leq b_i, 1 \leq i \leq k\}.$$

Write $\mathbf{x} = (x_1, \dots, x_k) \in \mathbb{R}^k$. Let

$$\delta = \left(\sum_{i=1}^k (b_i - a_i)^2 \right)^{1/2}$$

that is, the distance between the points (a_1, \dots, a_k) and (b_1, \dots, b_k) , which is the maximum distance between two points in I : for all $\mathbf{x}, \mathbf{y} \in I$,

$$|\mathbf{x} - \mathbf{y}| \leq \delta.$$

Suppose, for a contradiction, that I is not compact; that is, there exists an open cover $\mathcal{U} = \{U_i\}$ of I which contains no finite subcover of I .

For $1 \leq i \leq k$, let $c_i = \frac{1}{2}(a_i + b_i)$. The intervals $[a_i, c_i]$ and $[c_i, b_i]$ then determine 2^k k -cells Q_i whose union is I . At least one of these sets Q_i , call it I_1 , cannot be covered by any finite subcollection of \mathcal{U} (otherwise I could be so covered). We next subdivide I_1 and continue the process. We obtain a sequence (I_n) with the following properties:

- (i) $I \supset I_1 \supset I_2 \supset \dots$
- (ii) I_n is not covered by any finite subcollection of \mathcal{U}
- (iii) $|\mathbf{x} - \mathbf{y}| \leq 2^{-n}\delta$ for all $\mathbf{x}, \mathbf{y} \in I_n$

By (i) and 11.44, there is a point \mathbf{x}' which lies in every I_n . For some i , $\mathbf{x}' \in U_i$. Since U_i is open, there exists $r > 0$ such that $|\mathbf{y} - \mathbf{x}'| < r$ implies that $\mathbf{y} \in U_i$. If n is so large that $2^{-n}\delta < r$ (there is such an n , for otherwise $2^n \leq \frac{\delta}{r}$ for all positive integers n , which is absurd since \mathbb{R} is archimedean), then (iii) implies that $I_n \subset U_i$, which contradicts (ii). \square

We have now come to an important result, which will be crucial in proving the Heine–Borel theorem and Bolzano–Weierstrass theorem.

Proposition 11.47. *If $E \subset \mathbb{R}^k$ has one of the following three properties, then it has the other two:*

- (i) E is closed and bounded.
- (ii) E is compact.

(iii) Every infinite subset of E has a limit point in E .

Proof.

(i) \implies (ii) Suppose E is closed and bounded. Since E is bounded, then $E \subset I$ for some k -cell I . By 11.46, I is compact. Since E is a closed subset of a compact set, by 11.39, E is compact.

(ii) \implies (iii) This directly follows from 11.41.

(iii) \implies (i) If E is not bounded, then E contains points \mathbf{x}_n with

$$|\mathbf{x}_n| > n \quad (n = 1, 2, 3, \dots)$$

The set S consisting of these points \mathbf{x}_n is infinite and clearly has no limit point in \mathbb{R}^k , hence has none in E . Thus (iii) implies that E is bounded.

If E is not closed, then there is a point $\mathbf{x}_0 \in \mathbb{R}^k$ which is a limit point of E but not a point of E . For $n = 1, 2, 3, \dots$, there are points $\mathbf{x}_n \in E$ such that $|\mathbf{x}_n - \mathbf{x}_0| < \frac{1}{n}$. Let S be the set of these points \mathbf{x}_n . Then S is infinite (otherwise $|\mathbf{x}_n - \mathbf{x}_0|$ would have a constant positive value, for infinitely many n), S has \mathbf{x}_0 as a limit point, and S has no other limit point in \mathbb{R}^k . For if $\mathbf{y} \in \mathbb{R}^k$, $\mathbf{y} \neq \mathbf{x}_0$, then

$$\begin{aligned} |\mathbf{x}_n - \mathbf{y}| &\geq |\mathbf{x}_0 - \mathbf{y}| - |\mathbf{x}_n - \mathbf{x}_0| \\ &\geq |\mathbf{x}_0 - \mathbf{y}| - \frac{1}{n} \\ &\geq \frac{1}{2}|\mathbf{x}_0 - \mathbf{y}| \end{aligned}$$

for all but finitely many n ; this shows that \mathbf{y} is not a limit point of S (Theorem 2.20).

Thus S has no limit point in E ; hence E must be closed if (iii) holds. □

review
proof

Theorem 11.48 (Heine–Borel theorem). $E \subset \mathbb{R}^n$ is compact if and only if E is closed and bounded.

Proof. This is simply (i) \iff (ii) in the previous result. □

Bolzano–Weierstrass Theorem

Theorem 11.49 (Bolzano–Weierstrass theorem). Every bounded infinite subset of \mathbb{R}^n has a limit point in \mathbb{R}^n .

Proof. Suppose E is a bounded infinite subset of \mathbb{R}^n .

Since E is bounded, there exists an n -cell $I \subset \mathbb{R}^n$ such that $E \subset I$. Since I is compact, by 11.41, E has a limit point in I and thus \mathbb{R}^n . □

Cantor's Intersection Theorem

A collection \mathcal{A} of subsets of X is said to have the *finite intersection property* if the intersection of every finite subcollection of \mathcal{A} is non-empty.

Proposition 11.50. *Suppose $\mathcal{K} = \{K_i \mid i \in I\}$ is a collection of compact subsets of a metric space X , which satisfies the finite intersection property. Then $\bigcap_{i \in I} K_i \neq \emptyset$.*

Proof. We fix a member $K_1 \in \mathcal{K}$. Suppose, for a contradiction, that $\bigcap_{i \in I} K_i = \emptyset$; that is, no point of K_1 belongs to every $K_i \in \mathcal{K}$.

For $i \in I$, let $U_i = K_i^c$. Then the sets $\{U_i \mid i \in I\}$ form an open cover of K_1 . Since K_1 is compact by assumption, there exist finitely many indices i_1, \dots, i_n such that

$$K_1 \subset \bigcup_{k=1}^n U_{i_k}.$$

By de Morgan's laws, we have that

$$\bigcup_{k=1}^n U_{i_k} = \bigcup_{k=1}^n K_{i_k}^c = \left(\bigcap_{k=1}^n K_{i_k} \right)^c.$$

Thus

$$K_1 \subset \left(\bigcap_{k=1}^n K_{i_k} \right)^c,$$

which means that

$$K_1 \cap \bigcap_{k=1}^n K_{i_k} = \emptyset.$$

Thus $K_1, K_{i_1}, \dots, K_{i_n}$ is a finite subcollection of \mathcal{K} which has an empty intersection; this contradicts the finite intersection property of \mathcal{K} . □

Theorem 11.51 (Cantor's intersection theorem). *Suppose (K_n) is a decreasing sequence of non-empty compact sets; that is, $K_1 \supset K_2 \supset \dots$. Then $\bigcap_{n=1}^{\infty} K_n \neq \emptyset$.*

Proof. This is an immediate corollary of the previous result. □

The following result is a characterisation of compact sets.

Proposition 11.52. *K is compact if and only if every collection of closed subsets of K satisfies the finite intersection property.*

Proof.

\implies Suppose K is compact.

If \mathcal{U} is an open covering of K , then the collection \mathcal{F} of complements of sets in \mathcal{U} is a collection of closed sets whose intersection is empty (why?); and

conversely, if \mathcal{F} is a collection of closed sets whose intersection is empty, then the collection \mathcal{U} of complements of sets in \mathcal{F} is an open covering.

□

To complete proof

Sequential Compactness

Definition 11.53 (Sequential compactness). We say $K \subset X$ is *sequentially compact* if every sequence in K has a convergent subsequence in K .

We now show that compactness and sequential compactness are equivalent.

Proposition 11.54. $K \subset X$ is compact if and only if it is sequentially compact.

Proof.

\Rightarrow Suppose $K \subset X$ is compact. Take any sequence (y_n) from K . Suppose, for a contradiction, that every point $x \in K$ is not a limit of any subsequence of (y_n) . Then for all $x \in K$, there exists $r_x > 0$ such that $B_{r_x}(x)$ contains at most one point in (y_n) , which is x .

Consider the collection of open balls at each $x \in K$:

$$\{B_{r_x}(x) \mid x \in K\}.$$

This is an open cover of K . By the compactness of K , there exists a finite subcover of K :

$$\{B_{r_{x_1}}(x_1), \dots, B_{r_{x_N}}(x_N)\}.$$

In particular, these open balls cover $\{y_n\}$. Hence there must be some x_i ($1 \leq i \leq N$) such that there are infinitely many $y_j = x_i$. Consider the sequence (y_j) where each term in this sequence is equal to x_i ; this is a subsequence of (y_n) that converges to $x_i \in K$. This contradicts the assumption.

\Leftarrow Suppose, for a contradiction, that K is not compact. Then there exists an open cover $\{U_\alpha \mid \alpha \in \Lambda\}$ which has no finite subcover. Then Λ must be an infinite set.

If Λ is countable, WLOG, assume $\Lambda = \mathbb{N}$. Since any finite union

$$\bigcup_{i=1}^n U_i$$

cannot cover K , we can take some $x_n \in K \setminus \bigcup_{i=1}^n U_i$ for every $n \in \mathbb{N}$. Then we obtain a sequence (x_n) in K and so must have a convergent subsequence (x_{n_k}) that converges to some $x_0 \in K$. It follows that there must be some U_N such that $x_0 \in U_N$. Since U_N is open, there exists $r > 0$ such that

$$B_r(x_0) \subset U_N.$$

On the other hand, since $x_{n_k} \rightarrow x_0$, there exists $N' \in \mathbb{N}$ such that if $n_k \geq N'$ then

$$x_{n_k} \in B_r(x_0).$$

However, by our way of choosing x_n , whenever $n_k > \max\{N', N\}$, $x_{n_k} \notin U_N$. This leads to a contradiction. \square

§11.3 Perfect Sets

Definition and Uncountability

Definition 11.55 (Perfect set). E is *perfect* if

- (i) E is closed, and
- (ii) every point of E is a limit point of E .

Proposition 11.56. *Let non-empty $P \subset \mathbb{R}^k$ be perfect. Then P is uncountable.*

Proof. Since P has limit points, by 11.25, P is an infinite set.

Suppose, for a contradiction, that P is countable. This means we can list the points of P in a sequence:

$$\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots$$

Consider a sequence (B_n) of open balls, where B_n is any open ball centred at \mathbf{x}_n :

$$B_n = \{ \mathbf{y} \in \mathbb{R}^k \mid |\mathbf{y} - \mathbf{x}_n| < r \}.$$

Then its closure $\overline{B_n}$ is the closed ball

$$\overline{B_n} = \{ \mathbf{y} \in \mathbb{R}^k \mid |\mathbf{y} - \mathbf{x}_n| \leq r \}.$$

Suppose B_n has been constructed. Note that $B_n \cap P$ is not empty. Since P is perfect, every point of P is a limit point of P , so there exists B_{n+1} such that (i) $\overline{B_{n+1}} \subset B_n$, (ii) $\mathbf{x}_n \notin \overline{B_{n+1}}$, (iii) $B_{n+1} \cap P$ is not empty.

By (iii), B_{n+1} satisfies our induction hypothesis, and the construction can proceed.

Put $K_n = \overline{B_n} \cap P$. Since $\overline{B_n}$ is closed and bounded, $\overline{B_n}$ is compact. Since $\mathbf{x}_n \notin K_{n+1}$, no point of P lies in $\bigcap_{n=1}^{\infty} K_n$. Since $K_n \subset P$, this implies that $\bigcap_{n=1}^{\infty} K_n$ is empty. But each K_n is nonempty, by (iii), and $K_n \supset K_{n+1}$ by (i); this contradicts Cantor's intersection theorem (11.51). \square

Corollary. *Every interval $[a, b]$ is uncountable. In particular, \mathbb{R} is uncountable.*

Cantor Set

We now construct the Cantor set. Consider the interval

$$C_0 = [0, 1].$$

Remove the middle third $(\frac{1}{3}, \frac{2}{3})$ to give

$$C_1 = \left[0, \frac{1}{3}\right] \cup \left[\frac{2}{3}, 1\right].$$

Remove the middle thirds of these intervals to give

$$C_2 = \left[0, \frac{1}{9}\right] \cup \left[\frac{2}{9}, \frac{3}{9}\right] \cup \left[\frac{6}{9}, \frac{7}{9}\right] \cup \left[\frac{8}{9}, 1\right].$$

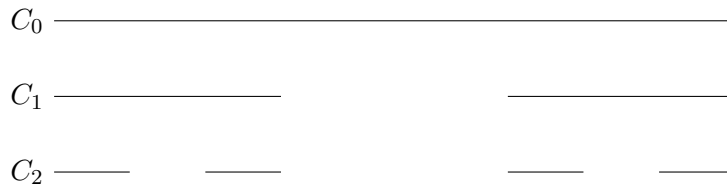


Figure 11.7: Cantor set

Repeating this process, we obtain a monotonically decreasing sequence of compact sets (C_n) , where C_n is the union of 2^n intervals, each of length 3^{-n} . Recursively, we have that $C_{n+1} = \frac{1}{3}C_n \cup (\frac{1}{3}C_n + \frac{2}{3})$.

Note that each C_n has the following properties:

- (i) closed (since each C_n is a finite union of closed sets, which is closed)
- (ii) compact (since each C_n is a closed subset of a compact set $[a, b]$)
- (iii) non-empty (since the endpoints 0 and 1 are in each C_n)

The **Cantor set** is defined to be the union

$$C := \bigcap_{n=1}^{\infty} C_n.$$

Lemma 11.57 (Properties of the Cantor set).

- (i) C is closed.
- (ii) C is compact.
- (iii) C is not empty.
- (iv) C has no interior points.

Proof.

- (i) C is the intersection of arbitrarily many closed sets, so C is closed.

- (ii) C is bounded in $[0, 1]$, by definition. Since C is closed and bounded, by the Heine–Borel theorem, C is compact.
- (iii) Since (C_n) is a decreasing sequence of non-empty compact sets, by Cantor's intersection theorem, $\bigcap_{n=1}^{\infty} C_n = C \neq \emptyset$.
- (iv) Suppose, for a contradiction, that there exists $p \in C$ which is an interior point. Then there exists some open interval around p , i.e., $p \in (a, b)$.

However in C_n , each interval has length $\frac{1}{3^n}$. Hence for any (a, b) we can find some $n \in \mathbb{N}$ such that (a, b) is not contained in C_n and hence not contained in C .

□

Proposition 11.58. C is a perfect set in \mathbb{R} which contains no open interval.

Proof. We will show that (i) C contains no open interval, and (ii) C is perfect.

- (i) No open interval of the form

$$\left(\frac{3k+1}{3^m}, \frac{3k+2}{3^m} \right),$$

where $k, m \in \mathbb{Z}^+$, has a point in common with C . Since every open interval (α, β) contains a open interval of the above form, if

$$3^{-m} < \frac{\beta - \alpha}{6},$$

C contains no open interval.

- (ii) Since we have shown that C is closed, it suffices to show that every point of C is a limit point.

Let $x \in C$, and let S be any open interval containing x . Let I_n be that interval of C_n which contains x . Choose n large enough, so that $I_n \subset S$. Let x_n be an endpoint of I_n , such that $x_n \neq x$.

It follows from the construction of C that $x_n \in C$. Hence x is a limit point of C , and C is perfect.

□

Corollary 11.59. C is uncountable.

One of the most interesting properties of the Cantor set is that it provides us with an example of an uncountable set of measure zero.

§11.4 Connectedness

Definition 11.60 (Connectedness). We say A and B are *separated* if

- (i) $A \cap \bar{B} = \emptyset$, and
- (ii) $\bar{A} \cap B = \emptyset$;

that is, no point of A lies in the closure of B , and no point of B lies in the closure of A . (Equivalently, no point of one set is a limit point of the other set.)

$E \subset X$ is *connected* if E is not the union of two non-empty separated sets.

Remark. Separated sets are of course disjoint, but disjoint sets need not be separated. For example, $[0, 1]$ and $(1, 2)$ are not separated, since 1 is a limit point of $(1, 2)$. However $(0, 1)$ and $(1, 2)$ are separated.

Example 11.61. In \mathbb{R}^2 , consider the set

$$E = \{(x, y) \mid x, y \in \mathbb{Q}\}.$$

Then E is not connected; if we let

$$A = \{(x, y) \mid x, y \in \mathbb{Q}, x < \sqrt{2}\},$$

$$B = \{(x, y) \mid x, y \in \mathbb{Q}, x > \sqrt{2}\},$$

then note that $A \cup B = E$, as well as $A \cap \bar{B} = \emptyset$ and $\bar{A} \cap B = \emptyset$.

Lemma 11.62. *Closed intervals in \mathbb{R} are connected.*

Proof. Suppose, for a contradiction, that a closed interval $[a, b]$ is not connected. Then there exists non-empty sets A and B , with $A \cap \bar{B} = \emptyset$ and $\bar{A} \cap B = \emptyset$. WLOG let $a \in A$.

Let $s = \sup A$. By 11.31, $s \in \bar{A}$. Then $\bar{A} \cap B = \emptyset$ implies $s \notin B$, so $s \in A$. Thus $A \cap \bar{B} = \emptyset$ implies $s \notin \bar{B}$. Hence there exists an open interval $(s - \varepsilon, s + \varepsilon)$ around s that is disjoint from B . But since $A \cup B = [a, b]$, we must have $(s - \varepsilon, s + \varepsilon) \subset A$. This contradicts the fact that s is the supremum of A . \square

The connected subsets of the real line have a particularly simple structure:

Lemma 11.63. *$E \subset \mathbb{R}$ is connected if and only if it has the following property: if $x, y \in E$ and $x < z < y$, then $z \in E$.*

Proof.

\Leftarrow If there exists $x, y \in E$ and some $z \in (x, y)$ such that $z \notin E$, then $E = A_z \cup B_z$ where

$$A_z = E \cap (-\infty, z), \quad B_z = E \cap (z, \infty).$$

Since $x \in A_z$ and $y \in B_z$, A and B are non-empty. Since $A_z \subset (-\infty, z)$ and $B_z \subset (z, \infty)$, they are separated. Hence E is not connected.

\Rightarrow Suppose E is not connected. Then there are non-empty separated sets A and B such that $A \cup B = E$. Pick $x \in A$, $y \in B$, and WLOG assume that $x < y$. Define

$$z := \sup(A \cap [x, y].)$$

By 11.31, $z \in \overline{A}$; hence $z \notin B$. In particular, $x \leq z < y$.

Case 1: $z \notin A$. It follows that $x < z < y$ and $z \notin E$.

Case 2: $z \in A$. Then $z \notin B$, hence there exists z_1 such that $z < z_1 < y$ and $z_1 \notin B$. Then $x < z_1 < y$ and $z_1 \notin E$.

□

Path Connectedness

§11.5 Separable Spaces

Definition 11.64 (Separable space). X is *separable* if it has a countable subset which is dense in X .

Example 11.65.

- \mathbb{R} is separable.

Proof. The set of rational numbers \mathbb{Q} is countable and is dense in \mathbb{R} . □

- \mathbb{C} is separable.

Proof. A countable dense subset of \mathbb{C} is the set of all complex numbers whose real and imaginary parts are both rational, i.e., the set $\{x + yi \mid x, y \in \mathbb{Q}\}$. □

- The *discrete metric space* X is separable if and only if X is countable.

Proof. The kind of metric implies that no proper subset of X can be dense in X . Hence the only dense set in X is X itself, and the statement follows. □

- The *sequence space* ℓ^∞ is the set of all bounded complex sequences, with the metric defined by

$$d(x, y) = \sup_{n \in \mathbb{N}} |x_n - y_n|.$$

ℓ^∞ is not separable.

§11.6 Baire Category Theorem

$E \subset X$ is called *nowhere dense* (in X) if the interior of the closure of A is empty, i.e., $(\overline{A})^\circ = \emptyset$.

Otherwise put, E is nowhere dense iff it is contained in a closed set with empty interior. Passing to complements, we can say equivalently that E is nowhere dense iff its complement contains a dense open set (why?).

Lemma 11.66. *Let X be a metric space.*

- (i) *Any subset of a nowhere dense set is nowhere dense.*
- (ii) *The union of finitely many nowhere dense sets is nowhere dense.*
- (iii) *The closure of a nowhere dense set is nowhere dense.*
- (iv) *If X has no isolated points, then every finite set is nowhere dense.*

Proof.

- (i)
- (ii)
- (iii)
- (iv)

□

Although the union of finitely many nowhere dense sets is nowhere dense, the union of countably many nowhere dense sets need not be nowhere dense: for instance, in $X = \mathbb{R}$, the rationals \mathbb{Q} are the union of countably many nowhere dense sets (why?), but the rationals are certainly not nowhere dense (indeed, they are everywhere dense, i.e. $(\overline{\mathbb{Q}})^\circ = \overline{\mathbb{Q}} = \mathbb{R}$).

This observation motivates the introduction of a larger class of sets: $A \subset X$ is called *meager* (or of first category) in X if it can be written as a countable union of nowhere dense sets; otherwise, it is *non-meager* (or of second category). The complement of a meager set is called *residual*.

We then have as an immediate consequence:

Lemma 11.67. *Let X be a metric space.*

- (i) *Any subset of a meager set is meager.*
- (ii) *The union of countably many meager sets is meager.*
- (iii) *If X has no isolated points, then every countable set is meager.*

We are now ready to state the Baire category theorem.

Theorem 11.68 (Baire category theorem). *Let X be a complete metric space.*

- (i) *A meager set has empty interior.*

(ii) *The complement of a meager set is dense. (That is, a residual set is dense.)*

(iii) *A countable intersection of dense open sets is dense.*

You should carefully verify that (i), (ii) and (iii) are equivalent statements, obtained by taking complements.

In applications we frequently need only the weak form of the Baire category theorem that is obtained by weakening “is dense” in (b,c) to “is non-empty” (which is valid whenever X is itself non-empty):

Corollary 11.69 (Weak form of the Baire category theorem). *Let X be a non-empty complete metric space.*

(i) *X cannot be written as a countable union of nowhere dense sets. (In other words, X is nonmeager in itself.)*

(ii) *If X is written as a countable union of closed sets, then at least one of those closed sets has nonempty interior.*

(iii) *A countable intersection of dense open sets is nonempty.*

Exercises

Exercise 11.1. Prove that the following are metrics.

(i) On an arbitrary set X , define

$$d(x, y) = \begin{cases} 1 & (x \neq y) \\ 0 & (x = y) \end{cases}$$

(This is called the *discrete metric*.)

(ii) On \mathbb{Z} , define $d(x, y)$ to be 2^{-m} , where 2^m is the largest power of two dividing $x - y$. The triangle inequality holds in the following stronger form, known as the ultrametric property:

$$d(x, z) \leq \max\{d(x, y), d(y, z)\}.$$

Indeed, this is just a rephrasing of the statement that if 2^m divides both $x - y$ and $y - z$, then 2^m divides $x - z$.

(This is called the *2-adic metric*. The role of 2 can be replaced by any other prime p , and the metric may also be extended in a natural way to the rationals \mathbb{Q} .)

(iii) Let $\mathcal{G} = (V, E)$ be a connected graph. Define d on V as follows: $d(v, v) = 0$, and $d(v, w)$ is the length of the shortest path from v to w .

(This is known as the *path metric*.)

(iv) Let G be a group generated by elements a, b and their inverses. Define a distance on G as follows: $d(v, w)$ is the minimal k such that $v = wg_1 \cdots g_k$, where $g_i \in \{a, b, a^{-1}, b^{-1}\}$ for all i .

(This is known as the *word metric*.)

(v) Let $X = \{0, 1\}^n$ (the boolean cube), the set of all strings of n zeroes and ones. Define $d(x, y)$ to be the number of coordinates in which x and y differ.

(This is known as the *Hamming distance*.)

(vi) Consider the set $P(\mathbb{R}^n)$ of one-dimensional subspaces of \mathbb{R}^n , that is to say lines through the origin. One way to define a distance on this set is to take, for lines L_1, L_2 , the distance between L_1 and L_2 to be

$$d(L_1, L_2) = \sqrt{1 - \frac{|\langle v, w \rangle|^2}{\|v\|^2 \|w\|^2}},$$

where v and w are any non-zero vectors in L_1 and L_2 respectively.

When $n = 2$, the distance between two lines is $\sin \theta$ where θ is the angle between those lines.

(This is known as the *projective space*.)

Exercise 11.2 (Product space). If (X, d_X) and (Y, d_Y) are metric spaces, set

$$d_{X \times Y}((x_1, y_1), (x_2, y_2)) = \sqrt{d_X(x_1, x_2)^2 + d_Y(y_1, y_2)^2}.$$

for $x_1, x_2 \in X, y_1, y_2 \in Y$.

Prove that $d_{X \times Y}$ gives a metric on $X \times Y$; we call $X \times Y$ the *product space*.

Solution. Reflexivity and symmetry are obvious. Less clear is the triangle inequality. We need to prove that

$$\begin{aligned} & \sqrt{d_X(x_1, x_3)^2 + d_Y(y_1, y_3)^2} + \sqrt{d_X(x_3, x_2)^2 + d_Y(y_3, y_2)^2} \\ & \geq \sqrt{d_X(x_1, x_2)^2 + d_Y(y_1, y_2)^2} \end{aligned} \quad (1)$$

Write $a_1 = d_X(x_2, x_3)$, $a_2 = d_X(x_1, x_3)$, $a_3 = d_X(x_1, x_2)$ and similarly $b_1 = d_Y(y_2, y_3)$, $b_2 = d_Y(y_1, y_3)$ and $b_3 = d_Y(y_1, y_2)$. Thus we want to show

$$\sqrt{a_2^2 + b_2^2} + \sqrt{a_1^2 + b_1^2} \geq \sqrt{a_3^2 + b_3^2}. \quad (2)$$

To prove this, note that from the triangle inequality we have $a_1 + a_2 \geq a_3$, $b_1 + b_2 \geq b_3$. Squaring and adding gives

$$a_1^2 + b_1^2 + a_2^2 + b_2^2 + 2(a_1a_2 + b_1b_2) \geq a_3^2 + b_3^2.$$

By Cauchy–Schwarz,

$$a_1a_2 + b_1b_2 \leq \sqrt{a_1^2 + b_1^2} \sqrt{a_2^2 + b_2^2}.$$

Substituting this into the previous line gives precisely the square of (2), and (1) follows. \square

12 Numerical Sequences and Series

Summary

- Sequences. Convergence, subsequences, Cauchy sequences. Limit superior and inferior.
- Series. Convergence tests.

Throughout, let (X, d) be a metric space.

§12.1 Sequences

Convergence

A **sequence** (a_n) in X is a function $f: \mathbb{N} \rightarrow X$ which maps $n \mapsto a_n$.

The *range* of a sequence (a_n) is the set

$$\{x \in X \mid \exists n \in \mathbb{N}, x = a_n\}.$$

Note that the range of a sequence may be a finite set or it may be infinite. (a_n) is *bounded* if its range is bounded.

Definition 12.1. A sequence (a_n) **converges** to $a \in X$, denoted by $a_n \rightarrow a$, if

$$\forall \varepsilon > 0, \quad \exists N \in \mathbb{N}, \quad \forall n \geq N, \quad d(a_n, a) < \varepsilon.$$

We call a a *limit* of (a_n) . If (a_n) does not converge, it is said to *diverge*.

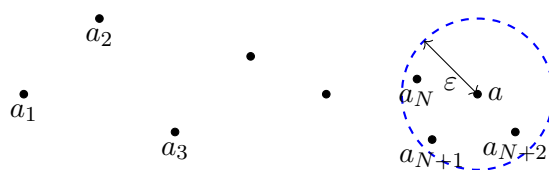


Figure 12.1: Convergence of sequence

Remark. This limit process conveys the intuitive idea that a_n can be made arbitrarily close to a , provided that n is sufficiently large. (Equivalently, if we remove more and more initial terms from the sequence, the *tail* of the sequence is increasingly closer to a .)

Remark. If $a_n \not\rightarrow a$, simply negate the definition for convergence:

$$\exists \varepsilon > 0, \quad \forall N \in \mathbb{N}, \quad \exists n \geq N, \quad d(a_n, a) \geq \varepsilon.$$

Remark. From the definition, the convergence of a sequence depends not only on the sequence itself, but also on the metric space X . For instance, the sequence given by $a_n = \frac{1}{n}$ converges in \mathbb{R} (to 0), but fails to converge in

\mathbb{R}^+ . In cases of possible ambiguity, we shall specify “convergent in X ” rather than “convergent”.

Example 12.2. $\frac{1}{n} \rightarrow 0$.

Proof. Fix $\varepsilon > 0$. By the Archimedean property, there exists $N \in \mathbb{N}$ such that $\frac{1}{N} < \varepsilon$. Take $N = \left\lfloor \frac{1}{\varepsilon} \right\rfloor + 1$. Then for all $n \geq N$,

$$\left| \frac{1}{n} - 0 \right| = \frac{1}{n} \leq \frac{1}{N} = \frac{1}{\left\lfloor \frac{1}{\varepsilon} \right\rfloor + 1} < \frac{1}{\frac{1}{\varepsilon}} = \varepsilon$$

as desired. Therefore $\frac{1}{n} \rightarrow 0$. □

A useful tip for finding the required N (in terms of ε) is to *work backwards* from the result we wish to show, as illustrated in the following example.

Example 12.3. Let $a_n = 1 + (-1)^n \frac{1}{\sqrt{n}}$. Then $a_n \rightarrow 1$.

Before our proof, we aim to find some $N \in \mathbb{N}$ such that if $n \geq N$ then

$$\begin{aligned} & |a_n - 1| < \varepsilon \\ \iff & \frac{1}{\sqrt{n}} = \left| (-1)^n \frac{1}{\sqrt{n}} \right| < \varepsilon \\ \iff & \frac{1}{n} < \varepsilon^2 \\ \iff & n > \frac{1}{\varepsilon^2} \end{aligned}$$

Hence take $N = \left\lfloor \frac{1}{\varepsilon^2} \right\rfloor + 1$.

Proof. Let $\varepsilon > 0$ be given. Take $N = \left\lfloor \frac{1}{\varepsilon^2} \right\rfloor + 1$. If $n \geq N$, then

$$\begin{aligned} |a_n - 1| &= \left| (-1)^n \frac{1}{\sqrt{n}} \right| = \frac{1}{\sqrt{n}} \\ &\leq \frac{1}{\sqrt{N}} = \frac{1}{\sqrt{\left\lfloor \frac{1}{\varepsilon^2} \right\rfloor + 1}} \\ &< \frac{1}{\sqrt{\frac{1}{\varepsilon^2}}} = \varepsilon \end{aligned}$$

as desired. Therefore $a_n \rightarrow 1$. □

Lemma 12.4 (Uniqueness of limit). *If a sequence converges, then its limit is unique.*

Proof. Let (a_n) be a sequence in X . Suppose that $a_n \rightarrow a$ and $a_n \rightarrow a'$ for $a, a' \in X$. We will show that $a' = a$.

Let $\varepsilon > 0$ be given. Then there exists $N, N' \in \mathbb{N}$ such that

$$\begin{aligned} n \geq N &\implies d(a_n, a) < \frac{\varepsilon}{2} \\ n \geq N' &\implies d(a_n, a') < \frac{\varepsilon}{2} \end{aligned}$$

Take $N_1 := \max\{N, N'\}$. If $n \geq N_1$, then both hold. By the triangle inequality,

$$d(a, a') \leq d(a, a_n) + d(a_n, a') < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon.$$

Since this holds for all $\varepsilon > 0$, we must have $d(a, a') = 0$. Hence $a = a'$. \square

Since the limit is unique, we can give it a notation.

Notation. If (a_n) converges to a , denote $\lim_{n \rightarrow \infty} a_n = a$.

We now outline some important properties of convergent sequences in metric spaces.

Lemma 12.5. *Let (a_n) be a sequence in X .*

- (i) $a_n \rightarrow a$ if and only if every open ball of a contains a_n for all but finitely many n .
- (ii) Every convergent sequence is bounded.
- (iii) Suppose $E \subset X$. Then a is a limit point of E if and only if there exists a sequence (a_n) in $E \setminus \{a\}$ such that $a_n \rightarrow a$.

Proof.

- (i) $\boxed{\implies}$ Suppose $a_n \rightarrow a$. Let $\varepsilon > 0$ be given, there exists $N \in \mathbb{N}$ such that

$$n \geq N \implies d(a_n, a) < \varepsilon \implies B_\varepsilon(a).$$

Hence $n \geq N$ implies $a_n \in B_\varepsilon(a)$.

$\boxed{\impliedby}$ Suppose every open ball of a contains all but finitely many of the a_n .

Let $\varepsilon > 0$ be given. Consider the open ball $B_\varepsilon(a)$. Since $B_\varepsilon(a)$ is an open ball of a , it will also eventually contain all a_n ; that is, there exists $N \in \mathbb{N}$ such that if $n \geq N$, then $a_n \in B_\varepsilon(a)$, i.e. $d(a_n, a) < \varepsilon$. Hence $a_n \rightarrow a$.

- (ii) Suppose $a_n \rightarrow a$. Take $\varepsilon = 1$, there exists $N \in \mathbb{N}$ such that

$$n \geq N \implies d(a_n, a) < 1.$$

Let

$$r = \max\{1, d(a_1, a), \dots, d(a_N, a)\},$$

then $d(a_n, a) \leq r$, so the range of a_n is bounded by $B_r(a)$. Hence (a_n) is bounded.

- (iii) $\boxed{\implies}$ Suppose a is a limit point of E .

Consider a sequence of open balls $(B_{\frac{1}{n}}(a))$, for $n \in \mathbb{N}$. Since a is a limit point, each open ball intersects with E at some point which is not a . We pick one such point a_n from each $B_{\frac{1}{n}}(a) \cap E$. Then

$$d(a_n, a) < \frac{1}{n}.$$

Let $\varepsilon > 0$ be given. By the Archimedean property, there exists $N \in \mathbb{N}$ such that $\frac{1}{N} < \varepsilon$. If $n \geq N$,

$$d(a_n, a) \leq \frac{1}{n} \leq \frac{1}{N} < \varepsilon,$$

which shows that $a_n \rightarrow a$.

\Leftarrow Suppose that there exists a sequence (a_n) in $E \setminus \{a\}$ such that $a_n \rightarrow a$. Then for each open ball $B_\varepsilon(a)$, we can find some $N \in \mathbb{N}$ such that if $n \in \mathbb{N}$ then

$$a_n \in B_\varepsilon(a).$$

Since $a_n \in E \setminus \{a\}$, this shows that a is a limit point of E .

□

Remark. A consequence of (ii) is its contrapositive: any unbounded sequence is divergent. Note that the converse is not true; a counterexample is $(-1)^n$.

Lemma 12.6 (Ordering). *Suppose (a_n) and (b_n) are convergent sequences, and $a_n \leq b_n$. Then*

$$\lim_{n \rightarrow \infty} a_n \leq \lim_{n \rightarrow \infty} b_n.$$

Proof. Let $a = \lim_{n \rightarrow \infty} a_n$, $b = \lim_{n \rightarrow \infty} b_n$. Suppose, for a contradiction, that $a > b$.

Let $\varepsilon = a - b > 0$ be given. There exists $N_1, N_2 \in \mathbb{N}$ such that

$$\begin{aligned} n \geq N_1 &\implies |a_n - a| < \frac{\varepsilon}{2}, \\ n \geq N_2 &\implies |b_n - b| < \frac{\varepsilon}{2}. \end{aligned}$$

Let $N = \max\{N_1, N_2\}$, then $n \geq N$ implies

$$a_n > a - \frac{\varepsilon}{2}, \quad b_n < b + \frac{\varepsilon}{2}$$

and thus

$$a_n - b_n > a - b - \varepsilon = 0$$

so $a_n > b_n$, which is a contradiction. □

Remark. If $a_n < b_n$, we may not necessarily have $\lim_{n \rightarrow \infty} a_n < \lim_{n \rightarrow \infty} b_n$. For instance, $-\frac{1}{n} < \frac{1}{n}$ but their limits are both 0.

Lemma 12.7 (Arithmetic properties). *Suppose (a_n) and (b_n) are convergent sequences in \mathbb{C} ; let*

$a = \lim_{n \rightarrow \infty} a_n$, $b = \lim_{n \rightarrow \infty} b_n$. Then

$$(i) \quad \lim_{n \rightarrow \infty} ca_n = ca, \text{ where } c \text{ is a constant} \quad (\text{scalar multiplication})$$

$$(ii) \quad \lim_{n \rightarrow \infty} (a_n + b_n) = a + b \quad (\text{addition})$$

$$(iii) \quad \lim_{n \rightarrow \infty} (a_n b_n) = ab \quad (\text{multiplication})$$

$$(iv) \quad \lim_{n \rightarrow \infty} \frac{a_n}{b_n} = \frac{a}{b} \quad (b_n \neq 0, b \neq 0) \quad (\text{division})$$

Proof.

- (i) The case where $c = 0$ is trivial. Now suppose $c \neq 0$. Let $\varepsilon > 0$ be given. Then there exists $N \in \mathbb{N}$ such that

$$n \geq N \implies |a_n - a| < \frac{\varepsilon}{|c|}.$$

Then if $n \geq N$,

$$|ca_n - ca| = |c| |a_n - a| < \varepsilon.$$

- (ii) Let $\varepsilon > 0$ be given. Since $a_n \rightarrow a$ and $b_n \rightarrow b$, there exists $N_1, N_2 \in \mathbb{N}$ such that

$$n \geq N_1 \implies |a_n - a| < \frac{\varepsilon}{2},$$

$$n \geq N_2 \implies |b_n - b| < \frac{\varepsilon}{2}.$$

Let $N = \max\{N_1, N_2\}$, then $n \geq N$ implies

$$\begin{aligned} |(a_n + b_n) - (a + b)| &\leq |a_n - a| + |b_n - b| \\ &< \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon. \end{aligned}$$

Hence $\lim_{n \rightarrow \infty} (a_n + b_n) = a + b$, as desired.

- (iii) Write

$$a_n b_n - ab = (a_n - a)(b_n - b) + a(b_n - b) + b(a_n - a).$$

Let $\varepsilon > 0$ be given. Since $a_n \rightarrow a$ and $b_n \rightarrow b$, there exist $N_1, N_2 \in \mathbb{N}$ such that

$$n \geq N_1 \implies |a_n - a| < \sqrt{\varepsilon},$$

$$n \geq N_2 \implies |b_n - b| < \sqrt{\varepsilon}.$$

Let $N = \max\{N_1, N_2\}$. Then $n \geq N$ implies

$$|(a_n - a)(b_n - b)| < \varepsilon,$$

and thus $\lim_{n \rightarrow \infty} (a_n - a)(b_n - b) = 0$.

Note that $\lim_{n \rightarrow \infty} a(b_n - b) = \lim_{n \rightarrow \infty} b(a_n - a) = 0$. Hence

$$\lim_{n \rightarrow \infty} (a_n b_n - ab) = 0.$$

- (iv) Since we have proven multiplication, it suffices to show that $\lim_{n \rightarrow \infty} \frac{1}{b_n} = \frac{1}{b}$.

Since $b_n \rightarrow b$, there exists $m \in \mathbb{N}$ such that

$$n \geq m \implies |b_n - b| < \frac{1}{2}|b|.$$

Let $\varepsilon > 0$ be given. There exists $N \in \mathbb{N}$, $N > m$ such that

$$n \geq N \implies |b_n - b| < \frac{1}{2}|b|^2 \varepsilon.$$

Hence for $n \geq N$,

$$\left| \frac{1}{b_n} - \frac{1}{b} \right| = \left| \frac{b - b_n}{b_n b} \right| < \frac{2}{|b|^2} |b_n - b| < \varepsilon.$$

□

We now prove the analogue for Euclidean spaces.

Lemma 12.8.

(i) Suppose $\mathbf{x}_n \in \mathbb{R}^k$ ($n = 1, 2, \dots$) and

$$\mathbf{x}_n = (\alpha_{1,n}, \dots, \alpha_{k,n}).$$

Then (\mathbf{x}_n) converges to $\mathbf{x} = (\alpha_1, \dots, \alpha_k)$ if and only if

$$\lim_{n \rightarrow \infty} \alpha_{i,n} = \alpha_i \quad (1 \leq i \leq k).$$

(ii) Suppose (\mathbf{x}_n) and (\mathbf{y}_n) are sequences in \mathbb{R}^k , (β_n) is a sequence of real numbers, and $\mathbf{x}_n \rightarrow \mathbf{x}$, $\mathbf{y}_n \rightarrow \mathbf{y}$, $\beta_n \rightarrow \beta$. Then

$$\lim_{n \rightarrow \infty} (\mathbf{x}_n + \mathbf{y}_n) = \mathbf{x} + \mathbf{y}, \quad \lim_{n \rightarrow \infty} \mathbf{x}_n \cdot \mathbf{y}_n = \mathbf{x} \cdot \mathbf{y}, \quad \lim_{n \rightarrow \infty} \beta_n \mathbf{x}_n = \beta \mathbf{x}.$$

Proof.

(i) \Rightarrow Suppose $\mathbf{x}_n \rightarrow \mathbf{x}$. From the definition of the norm in \mathbb{R}^k , the inequalities

$$|\alpha_{i,n} - \alpha_i| \leq \|\mathbf{x}_n - \mathbf{x}\|$$

follow immediately, which show that

$$\lim_{n \rightarrow \infty} \alpha_{i,n} = \alpha_i \quad (1 \leq i \leq k).$$

\Leftarrow Suppose $\lim_{n \rightarrow \infty} \alpha_{i,n} = \alpha_i$ for $i = 1, \dots, k$. Then to each $\varepsilon > 0$, there exists $N \in \mathbb{N}$ such that $n \geq N$ implies

$$|\alpha_{i,n} - \alpha_i| < \frac{\varepsilon}{\sqrt{k}} \quad (i = 1, \dots, k).$$

Hence $n \geq N$ implies

$$\|\mathbf{x}_n - \mathbf{x}\| = \left(\sum_{i=1}^k |\alpha_{i,n} - \alpha_i|^2 \right)^{1/2} < \varepsilon,$$

so that $\mathbf{x}_n \rightarrow \mathbf{x}$.

(ii) This follows from (i) and 12.7.

□

The next result provides a useful method to evaluate limits of sequences.

Lemma 12.9 (Squeeze theorem). *Let $a_n \leq c_n \leq b_n$ where (a_n) and (b_n) are convergent sequences*

such that $\lim_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} b_n = L$. Then (c_n) is also a convergent sequence, and

$$\lim_{n \rightarrow \infty} c_n = L.$$

Proof. Let $\varepsilon > 0$ be given. There exist $N_1, N_2 \in \mathbb{N}$ such that

$$n \geq N_1 \implies |a_n - L| < \varepsilon,$$

$$n \geq N_2 \implies |b_n - L| < \varepsilon.$$

In particular, we have

$$a_n > L - \varepsilon, \quad b_n < L + \varepsilon.$$

Let $N = \max\{N_1, N_2\}$. Then $n \geq N$ implies

$$L - \varepsilon < a_n \leq c_n \leq b_n < L + \varepsilon$$

or

$$|c_n - L| < \varepsilon.$$

Hence (c_n) is convergent, and $c_n \rightarrow L$. □

The following example is a classic application of the squeeze theorem.

Example 12.10. Show that $\lim_{n \rightarrow \infty} \frac{\sin n}{n} = 0$.

Proof. We have $-1 \leq \sin n \leq 1$, so

$$-\frac{1}{n} \leq \frac{\sin n}{n} \leq \frac{1}{n}.$$

Now

$$\lim_{n \rightarrow \infty} \frac{1}{n} = \lim_{n \rightarrow \infty} \left(-\frac{1}{n}\right) = 0,$$

so the squeeze theorem yields the desired result. □

Subsequences

Definition 12.11 (Subsequence). Given a sequence (a_n) , consider a sequence (n_k) of positive integers such that $n_1 < n_2 < \dots$. Then (a_{n_k}) is called a **subsequence** of (a_n) . If (a_{n_k}) converges, its limit is called a *subsequential limit* of (a_n) .

Proposition 12.12. (a_n) converges to a if and only if every subsequence of (a_n) converges to a .

Proof.

\Rightarrow Suppose $a_n \rightarrow a$. Let $\varepsilon > 0$ be given. Then there exists $N \in \mathbb{N}$ such that

$$n \geq N \implies d(a_n, a) < \varepsilon.$$

Every subsequence of (a_n) can be written in the form (a_{n_k}) where $n_1 < n_2 < \dots$ is a strictly increasing sequence of positive integers. Pick M such that $n_M \geq N$. Then

$$k > M \implies n_k > n_M \geq N \implies d(a_{n_k}, a) < \varepsilon.$$

Hence every subsequence of (a_n) converges to a .

\Leftarrow Suppose every subsequence of (a_n) converges to a . Since (a_n) is a subsequence of itself, we must have $a_n \rightarrow a$. □

Proposition 12.13. In a compact metric space, any sequence has a convergent subsequence.

Proof. Suppose (a_n) is a sequence in a compact metric space X .

Let E be the range of (a_n) . We consider two cases:

Case 1: E is finite. Notice that there are infinitely many terms in the sequence (a_n) , but only finitely many distinct terms in E . By the pigeonhole principle, at least one term of E appears infinitely many times in the sequence.

That is, there exists $a \in E$ and a sequence (n_k) with $n_1 < n_2 < \dots$ such that

$$a_{n_1} = a_{n_2} = \dots = a.$$

This subsequence (a_{n_k}) evidently converges to a .

Case 2: E is infinite. If E is infinite, then E is an infinite subset of a compact set. By 11.41, E has a limit point $a \in X$.

We now construct a subsequence (a_{n_k}) of (a_n) such that $a_{n_k} \rightarrow a$.

- Choose n_1 so that $d(a, a_{n_1}) < 1$.
- Having chosen n_1, \dots, n_{k-1} , choose n_k where $n_k > n_{k-1}$ such that $d(a, a_{n_k}) < \frac{1}{k}$ (such n_k exists due to 11.25).

Then $a_{n_k} \rightarrow a$.

□

Corollary 12.14 (Bolzano–Weierstrass). *Every bounded sequence in \mathbb{R}^k has a convergent subsequence.*

Proof. By 11.47, every bounded sequence in \mathbb{R}^k lives in a compact subset of \mathbb{R}^k , and therefore it lives in a compact metric space. Hence by the previous result, it contains a convergent subsequence converging to a point in \mathbb{R}^k . □

Lemma 12.15. *Suppose (a_n) is a sequence in X . Then the subsequential limits of (a_n) form a closed subset of X .*

Proof. Let E be the set of all subsequential limits of (a_n) , let q be a limit point of E . We want to show that $q \in E$.

Choose n_1 so that $a_{n_1} \neq q$. (If no such n_1 exists, then E has only one point, and there is nothing to prove.) Put $\delta = d(q, a_{n_1})$. Suppose n_1, \dots, n_{i-1} are chosen. Since q is a limit point of E , there is an $a \in E$ with $d(a, q) < 2^{-1}\delta$. Since $a \in E$, there is an $n_i > n_{i-1}$ such that $d(a, a_{n_i}) < 2^{-i}\delta$. Thus

$$d(q, a_{n_i}) < 2^{1-i}\delta$$

for $i = 1, 2, 3, \dots$. This says that (a_{n_i}) converges to q . Hence $q \in E$. □

Cauchy Sequences

This is a very helpful way to determine whether a sequence is convergent or divergent, as it does not require the limit to be known. Subsequently we will see many instances where the convergence of all sorts of limits are compared with similar counterparts; generally we describe such properties as *Cauchy criteria*.

Definition 12.16 (Cauchy sequence). A sequence (a_n) in X is a **Cauchy sequence** if

$$\forall \varepsilon > 0, \quad \exists N \in \mathbb{N}, \quad \forall n, m \geq N, \quad d(a_n, a_m) < \varepsilon.$$

Remark. Intuitively, the distances between any two terms becomes sufficiently small after a certain point.

A natural question is regarding the relationship between convergent sequences and Cauchy sequences. We now address this.

Proposition 12.17.

- (i) In any metric space, every convergent sequence is a Cauchy sequence.
- (ii) If X is a compact metric space and if (a_n) is a Cauchy sequence in X , then (a_n) converges to some point of X .
- (iii) In \mathbb{R}^k , every Cauchy sequence converges.

Remark. The converse of (i) is not true. For instance, the sequence $\{3, 3.1, 3.14, 3.141, 3.1415, \dots\}$ is a Cauchy sequence but does not converge in \mathbb{Q} .

Proof.

- (i) Suppose $a_n \rightarrow a$. Let $\varepsilon > 0$. There exists $N \in \mathbb{N}$ such that for all $n \geq N$,

$$d(a_n, a) < \frac{\varepsilon}{2}.$$

Then for all $n, m \geq N$,

$$d(a_n, a_m) \leq d(a_n, a) + d(a_m, a) < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon,$$

as desired. Hence (a_n) is a Cauchy sequence.

- (ii) Let (a_n) be a Cauchy sequence in X . Since X is compact, it is sequentially compact. Then there exists a subsequence (a_{n_k}) such that $a_{n_k} \rightarrow a$.

Claim. $a_n \rightarrow a$.

Let $\varepsilon > 0$. Since (a_n) is a Cauchy sequence, there exists $N_1 \in \mathbb{N}$ such that

$$n, m \geq N_1 \implies d(a_n - a_m) < \frac{\varepsilon}{2}.$$

$a_{n_k} \rightarrow a$ implies there exists $N_2 \in \mathbb{N}$ such that

$$n_k \geq N_2 \implies d(a_{n_k}, a) < \frac{\varepsilon}{2}.$$

Let $N = \max\{N_1, N_2\}$, fix some $n_k \geq N$. Then $n \geq N$ implies

$$d(a_n, a) \leq d(a_n, a_{n_k}) + d(a_{n_k}, a) < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon.$$

(iii) Suppose (a_n) is a Cauchy sequence.

We perform three steps:

- We first show that (a_n) is bounded:

Pick $N \in \mathbb{N}$ such that $|a_n - a_N| \leq 1$ for all $n \geq N$. Then

$$|a_n| \leq \max\{1 + |a_N|, |a_1|, \dots, |a_{N-1}|\}.$$

- Since (a_n) is bounded, by Bolzano–Weierstrass, (a_n) contains a subsequence (a_{n_k}) which converges to a .
- We now show that $a_n \rightarrow a$.

Let $\varepsilon > 0$ be given. Since (a_n) is a Cauchy sequence, there exists $N_1 \in \mathbb{N}$ such that

$$n, m \geq N_1 \implies |a_n - a_m| < \frac{\varepsilon}{2}.$$

Since $a_{n_k} \rightarrow a$, there exists $M \in \mathbb{N}$ such that for all $k > M$,

$$n_k > n_M \implies |a_{n_k} - a| < \frac{\varepsilon}{2}.$$

Now since $n_1 < n_2 < \dots$ is a sequence of strictly increasing positive integers, we can pick $i > M$ such that $n_k > N_1$. Then for all $n \geq N_1$, by setting $m = n_k$ we obtain

$$|a_n - a_{n_k}| < \frac{\varepsilon}{2}, \quad |a_{n_k} - a| < \frac{\varepsilon}{2}.$$

Hence

$$|a_n - a| \leq |a_n - a_{n_k}| + |a_{n_k} - a| < \varepsilon.$$

Therefore (a_n) is convergent, and $a_n \rightarrow a$.

□

Definition 12.18. A metric space X is *complete* if every Cauchy sequence in X converges.

Remark. The above result shows that all compact metric spaces and all Euclidean spaces are complete. It also implies that every closed subset E of a complete metric space X is complete. (Every Cauchy sequence in E is a Cauchy sequence in X , hence it converges to some $a \in X$, and actually $a \in E$ since E is closed.)

Example 12.19. The sequence (a_n) is defined as follows:

$$a_n = 1 + \frac{1}{2} + \dots + \frac{1}{n}.$$

(a_n) does not converge in \mathbb{R} .

Proof. We claim that (a_n) is not a Cauchy sequence. WLOG assume $n > m$. Consider

$$|a_n - a_m| = \frac{1}{m+1} + \frac{1}{m+2} + \dots + \frac{1}{n} \geq \frac{n-m}{n} = 1 - \frac{m}{n}.$$

Let $n = 2m$, then

$$|a_n - a_m| = |a_{2m} - a_m| > \frac{1}{2}.$$

Hence (a_n) is not a Cauchy sequence, so it does not converge. □

Monotonic Sequences

Definition 12.20 (Monotonic sequence). A sequence (a_n) in \mathbb{R} is

- (i) *monotonically increasing* if $a_n \leq a_{n+1}$ for $n \in \mathbb{N}$;
- (ii) *monotonically decreasing* if $a_n \geq a_{n+1}$ for $n \in \mathbb{N}$;
- (iii) **monotonic** if it is either monotonically increasing or monotonically decreasing.

Lemma 12.21 (Monotone convergence theorem). *A monotonic sequence in \mathbb{R} converges if and only if it is bounded.*

Proof. We show the case for monotonically increasing sequences; the case for monotonically decreasing sequences is similar.

\Rightarrow We already proved that a convergent sequence is bounded.

\Leftarrow Suppose (a_n) is a monotonically increasing sequence bounded above.

Let E be the range of a_n . Since E is bounded above, let $a = \sup E$.

Claim. $a_n \rightarrow a$.

By definition of supremum, $a_n \leq a$ for all $n \in \mathbb{N}$. For every $\varepsilon > 0$, there exists $N \in \mathbb{N}$ such that

$$a - \varepsilon < a_N \leq a,$$

otherwise $a - \varepsilon$ would be an upper bound of E . Since (a_n) is monotonically increasing, $n \geq N$ implies $a_N \leq a_n \leq a$, so

$$a - \varepsilon < a_n \leq a,$$

which implies $|a_n - a| < \varepsilon$. Hence $a_n \rightarrow a$. □

Limit Superior and Inferior

For properly divergent sequences, we make the following definition.

Definition 12.22. Suppose (a_n) is a sequence in \mathbb{R} . We write $a_n \rightarrow \infty$ if

$$\forall M \in \mathbb{R}, \quad \exists N \in \mathbb{N}, \quad \forall n \geq N, \quad a_n \geq M.$$

Similarly, we write $a_n \rightarrow -\infty$ if

$$\forall M \in \mathbb{R}, \quad \exists N \in \mathbb{N}, \quad \forall n \geq N, \quad a_n \leq M.$$

Definition 12.23. Suppose (a_n) is a sequence in $[-\infty, \infty]$. Define respectively the *limit superior* and *limit infimum* of (a_n) as

$$\begin{aligned} \limsup_{n \rightarrow \infty} a_n &:= \lim_{n \rightarrow \infty} \sup\{a_n, a_{n+1}, \dots\} \\ \liminf_{n \rightarrow \infty} a_n &:= \lim_{n \rightarrow \infty} \inf\{a_n, a_{n+1}, \dots\} \end{aligned}$$

Remark. The limit superior and limit infimum exist due to the existence of supremum and infimum in $\overline{\mathbb{R}}$.

Remark. [Rud76] defines the limit superior and infimum in another manner, using subsequential limits; both definitions are equivalent.

Example 12.24.

- Let (a_n) be a sequence containing all rationals. Then every real number is a subsequential limit, and

$$\limsup_{n \rightarrow \infty} a_n = +\infty, \quad \liminf_{n \rightarrow \infty} a_n = -\infty.$$

- Let $a_n = \frac{(-1)^n}{1 + \frac{1}{n}}$. Then

$$\limsup_{n \rightarrow \infty} a_n = 1, \quad \liminf_{n \rightarrow \infty} a_n = -1.$$

Lemma 12.25.

$$\liminf_{n \rightarrow \infty} a_n = -\limsup_{n \rightarrow \infty} (-a_n).$$

Proof. Exercise; use the definitions and 10.12. □

Lemma 12.26. A sequence (a_n) in $[-\infty, \infty]$ converges if and only if

$$\limsup_{n \rightarrow \infty} a_n = \liminf_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} a_n.$$

Proof. □

Proposition 12.27. Suppose (a_n) is a sequence in \mathbb{R} . Then

$$(i) \limsup_{n \rightarrow \infty} a_n \in E;$$

(ii) if $a > \limsup_{n \rightarrow \infty} a_n$, there exists $N \in \mathbb{N}$ such that $a_n < a$ for all $n \geq N$.

Moreover, $\limsup_{n \rightarrow \infty} a_n$ is the only number that satisfies (i) and (ii).

Proof.

(i) We consider three cases for the value of $\limsup_{n \rightarrow \infty} a_n$:

- If $\limsup_{n \rightarrow \infty} a_n = +\infty$, then $\sup E = +\infty$, so E is not bounded above. Hence (a_n) is not bounded above, so (a_n) has a subsequence (a_{n_k}) such that $a_{n_k} \rightarrow \infty$
- If $\limsup_{n \rightarrow \infty} a_n \in \mathbb{R}$, then $\sup E \in \mathbb{R}$, so E is bounded above. Hence at least one subsequential limit exists, so that (i) follows from Theorems 3.7 and 2.28.
- If $\limsup_{n \rightarrow \infty} a_n = -\infty$, then $\sup E = -\infty$, so E contains only one element, namely $-\infty$. Hence (a_n) has no subsequential limit. Thus for any $M \in \mathbb{R}$, $a_n > M$ for at most a finite number of values of n , so that $a_n \rightarrow -\infty$.

(ii) We prove by contradiction.

Suppose there is a number $a > \limsup_{n \rightarrow \infty} a_n$ such that $a_n \geq a$ for infinitely many values of n . In that case, there is a number $y \in E$ such that $y \geq a > \limsup_{n \rightarrow \infty} a_n$, contradicting the definition of $\limsup_{n \rightarrow \infty} a_n$.

We now show uniqueness. Suppose, for a contradiction, that two numbers p and q satisfy (i) and (ii). WLOG assume $p < q$. Then choose a such that $p < a < q$. Since p satisfies (i), we have $a_n < a$ for all $n \geq N$. But then q cannot satisfy (i). \square

Of course, an analogous result is true for $\liminf_{n \rightarrow \infty} a_n$.

Lemma 12.28 (Comparison). If $a_n \leq b_n$ for $n \geq N$ (where N is fixed), then

$$\liminf_{n \rightarrow \infty} a_n \leq \liminf_{n \rightarrow \infty} b_n,$$

$$\limsup_{n \rightarrow \infty} a_n \leq \limsup_{n \rightarrow \infty} b_n.$$

Lemma 12.29 (Arithmetic properties).

$$(i) \text{ If } k > 0, \limsup_{n \rightarrow \infty} ka_n = k \limsup_{n \rightarrow \infty} a_n.$$

$$\text{If } k < 0, \limsup_{n \rightarrow \infty} ka_n = k \liminf_{n \rightarrow \infty} a_n.$$

$$(ii) \limsup_{n \rightarrow \infty} (a_n + b_n) \leq \limsup_{n \rightarrow \infty} a_n + \limsup_{n \rightarrow \infty} b_n$$

Moreover, $\limsup_{n \rightarrow \infty} (a_n + b_n)$ may be bounded from below as follows:

$$\limsup_{n \rightarrow \infty} (a_n + b_n) \geq \limsup_{n \rightarrow \infty} a_n + \liminf_{n \rightarrow \infty} b_n.$$

write down the analogous properties for \liminf , and to prove (i) and (ii)

Now you should try to prove (i) for \liminf as well; as for (ii), try to explain why properties (i),(ii) for \limsup and property (i) for \liminf would imply property (ii) for \liminf

§12.2 Series

Definition 12.30 (Series). Given a sequence (a_n) , we associate a sequence (s_n) , where

$$s_n = \sum_{k=1}^n a_k = a_1 + a_2 + \cdots + a_n,$$

where the term s_n is called the n -th *partial sum*. The sequence (s_n) is often written as

$$\sum_{n=1}^{\infty} a_n,$$

which we call a *series*.

Definition 12.31 (Convergence of series). We say that the series *converges* if $s_n \rightarrow s$ (the sequence of partial sums converges), and write $\sum_{n=1}^{\infty} a_n = s$; that is,

$$\forall \varepsilon > 0, \quad \exists N \in \mathbb{N}, \quad \forall n \geq N, \quad \left| \sum_{k=1}^n a_k - s \right| < \varepsilon.$$

The number s is called the *sum* of the series. If (s_n) diverges, the series is said to *diverge*.

Notation. When there is no possible ambiguity, we write $\sum_{n=1}^{\infty} a_n$ simply as $\sum a_n$.

The Cauchy criterion can be restated in the following form:

Lemma 12.32 (Cauchy criterion). $\sum a_n$ converges if and only if

$$\forall \varepsilon > 0, \quad \exists N \in \mathbb{N}, \quad \forall n \geq m \geq N, \quad \left| \sum_{k=m}^n a_k \right| \leq \varepsilon.$$

Convergence Tests

To determine the convergence of a series, apart from using the definition and the Cauchy criterion, we also have the following methods:

- Divergence test (12.33)
- Boundedness of partial sums (12.34, for series of non-negative terms)
- Comparison test (12.35)
- Root test (12.39)
- Ratio test (12.40)
- Absolute convergence (12.41)

Lemma 12.33 (Divergence test). *If $a_n \not\rightarrow 0$, then $\sum a_n$ diverges.*

Proof. We prove the contrapositive: if $\sum a_n$ converges, then $a_n \rightarrow 0$.

In the Cauchy criterion, take $m = n$, then $|a_n| \leq \varepsilon$ for all $n \geq N$. □

Remark. The converse is not true; a counterexample of the harmonic series.

Lemma 12.34. *A series of non-negative terms converges if and only if its partial sums form a bounded sequence.*

Proof. Partial sums are monotonically increasing. But bounded monotonic sequences converge. □

Lemma 12.35 (Comparison test). *Consider two sequences (a_n) and (b_n) .*

(i) *Suppose $|a_n| \leq b_n$ for all $n \geq N_0$ (where N_0 is some fixed integer). If $\sum b_n$ converges, then $\sum a_n$ converges.*

(ii) *Suppose $a_n \geq b_n \geq 0$ for all $n \geq N_0$. If $\sum b_n$ diverges, then $\sum a_n$ diverges.*

Proof.

(i) Since $\sum b_n$ converges, by the Cauchy criterion, fix $\varepsilon > 0$, there exists $N \in \mathbb{N}$, $N \geq N_0$ such that for $n \geq m \geq N$,

$$\sum_{k=m}^n b_k \leq \varepsilon.$$

By the triangle inequality,

$$\left| \sum_{k=m}^n a_k \right| \leq \sum_{k=m}^n |a_k| \leq \sum_{k=m}^n b_k \leq \varepsilon,$$

so $\sum a_n$ converges, by the Cauchy criterion.

(ii) We prove the contrapositive. If $\sum a_n$ converges, and since $|b_n| \leq a_n$ for all $n \geq N_0$, then by (i), $\sum b_n$ converges.

□

To employ the comparison test, we need to be familiar with several series whose convergence or divergence is known.

Example 12.36 (Geometric series). A geometric series takes the form

$$\sum_{n=0}^{\infty} x^n.$$

Proposition.

(i) If $|x| < 1$, then $\sum x^n$ converges;

$$\sum_{n=0}^{\infty} x^n = \frac{1}{1-x}.$$

(ii) If $|x| \geq 1$, then $\sum x^n$ diverges.

Proof.

(i) For $|x| < 1$, the n -th partial sum is given by

$$\sum_{k=0}^n x^k = 1 + x + x^2 + \cdots + x^n. \quad (1)$$

Multiplying both sides of (1) by x gives

$$x \sum_{k=0}^n x^k = x + x^2 + x^3 \cdots + x^{n+1}. \quad (2)$$

Taking the difference of (1) and (2),

$$(1-x) \sum_{k=0}^n x^k = 1 - x^{n+1}$$

and so

$$\sum_{k=0}^n x^k = \frac{1 - x^{n+1}}{1-x}.$$

Taking limits $n \rightarrow \infty$, the result follows.

(ii) For $|x| \geq 1$, $x^n \not\rightarrow 0$. By the divergence test, $\sum x^n$ diverges.

□

Example 12.37 (p -series). A p -series takes the form

$$\sum_{n=1}^{\infty} \frac{1}{n^p}.$$

To determine the convergence of p -series, we first prove the following lemma, which states that a rather “thin” subsequence of (a_n) determines the convergence of $\sum a_n$.

Lemma (Cauchy condensation test). *Suppose $a_1 \geq a_2 \geq \dots \geq 0$. Then $\sum a_n$ converges if and only if the series*

$$\sum_{k=0}^{\infty} 2^k a_{2^k} = a_1 + 2a_2 + 4a_4 + \dots$$

converges.

Proof. Let s_n and t_k denote the n -th partial sum of (a_n) and the k -th partial sum of $(2^k a_{2^k})$ respectively; that is,

$$\begin{aligned} s_n &= a_1 + a_2 + \dots + a_n, \\ t_k &= a_1 + 2a_2 + \dots + 2^k a_{2^k}. \end{aligned}$$

We consider two cases:

- For $n < 2^k$, group terms to give

$$\begin{aligned} s_n &= a_1 + a_2 + \dots + a_n \\ &\leq a_1 + (a_2 + a_3) + \dots + (a_{2^k} + \dots + a_{2^{k+1}-1}) \\ &\leq a_1 + 2a_2 + \dots + 2^k a_{2^k} \\ &= t_k. \end{aligned}$$

By comparison test, if (t_k) converges, then (s_n) converges.

- For $n > 2^k$,

$$\begin{aligned} s_n &\geq a_1 + a_2 + (a_3 + a_4) + \dots + (a_{2^{k-1}+1} + \dots + a_{2^k}) \\ &\geq \frac{1}{2} a_1 + a_2 + 2a_4 + \dots + 2^{k-1} a_{2^k} \\ &= \frac{1}{2} t_k. \end{aligned}$$

By comparison test, if (s_n) converges, then (t_k) converges.

□

Proposition (p -test).

(i) If $p > 1$, $\sum \frac{1}{n^p}$ converges.

(ii) If $p \leq 1$, $\sum \frac{1}{n^p}$ diverges.

Proof. Note that if $p \leq 0$, then $\frac{1}{n^p} \not\rightarrow 0$. By the divergence test, $\sum \frac{1}{n^p}$ diverges.

If $p > 0$, we want to apply the above lemma. Consider the series

$$\sum_{k=0}^{\infty} 2^k \cdot \frac{1}{(2^k)^p} = \sum_{k=0}^{\infty} 2^{(1-p)k} = \sum_{k=0}^{\infty} (2^{1-p})^k,$$

which is a geometric series. Hence the above series converges if and only if $|2^{1-p}| < 1$, which holds if and only if $1 - p < 0$. Then apply the above lemma to conclude the convergence of $\frac{1}{n^p}$. □

Remark. If $p = 1$, the resulting series is known as the *harmonic series* (which diverges). If $p = 2$, the resulting series converges, and the sum of this series is $\frac{\pi^2}{6}$ (Basel problem).

Example 12.38 (The number e). Consider the series

$$\sum_{n=0}^{\infty} \frac{1}{n!}.$$

We first show that the above series converges. Consider the n -th partial sum:

$$\begin{aligned} \sum_{k=0}^n \frac{1}{k!} &= \frac{1}{0!} + \frac{1}{1!} + \frac{1}{2!} + \frac{1}{3!} + \cdots + \frac{1}{n!} \\ &\leq 1 + 1 + \frac{1}{2} + \frac{1}{2^2} + \cdots + \frac{1}{2^{n-1}} \\ &< 1 + 1 + \frac{1}{2} + \frac{1}{2^2} + \cdots = 3. \end{aligned}$$

Since the partial sums are bounded (by 3), and the terms are non-negative, the series converges. Then we can make the following definition for the sum of the series:

$$e := \sum_{n=0}^{\infty} \frac{1}{n!}$$

Proposition. e is irrational.

Proof. Suppose, for a contradiction, that e is rational. Then $e = \frac{p}{q}$, where p and q are positive integers. Let s_n denote the n -th partial sum:

$$s_n = \sum_{k=0}^n \frac{1}{k!}.$$

Then

$$\begin{aligned} e - s_n &= \frac{1}{(n+1)!} + \frac{1}{(n+2)!} + \frac{1}{(n+3)!} + \cdots \\ &< \frac{1}{(n+1)!} \left(1 + \frac{1}{n+1} + \frac{1}{(n+1)^2} + \cdots \right) \\ &= \frac{1}{(n+1)!} \cdot \frac{n+1}{n} = \frac{1}{n!n} \end{aligned}$$

and thus

$$0 < e - s_n < \frac{1}{n!n}.$$

Taking $n = q$ and multiplying both sides by $q!$ gives

$$0 < q!(e - s_q) < \frac{1}{q}.$$

Note that $q!e$ is an integer (by assumption), and

$$q!s_q = q! \left(1 + 1 + \frac{1}{2!} + \cdots + \frac{1}{q!} \right)$$

is an integer, so $q!(e - s_q)$ is an integer. Since $q \geq 1$, this implies the existence of an integer between 0 and 1, which is absurd. Hence we have reached a contradiction. \square

Lemma. e is equivalent to the following:

$$\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n = e.$$

Proof. Let

$$s_n = \sum_{k=0}^n \frac{1}{k!}, \quad t_n = \left(1 + \frac{1}{n}\right)^n.$$

By the binomial theorem,

$$t_n = 1 + 1 + \frac{1}{2!} \left(1 - \frac{1}{n}\right) + \frac{1}{3!} \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) + \cdots + \frac{1}{n!} \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{n-1}{n}\right).$$

Comparing term by term, we see that $t_n \leq s_n$. By 12.28, we have that

$$\limsup_{n \rightarrow \infty} t_n \leq \limsup_{n \rightarrow \infty} s_n = e.$$

Next, if $n \geq m$,

$$t_n \geq 1 + 1 + \frac{1}{2!} \left(1 - \frac{1}{n}\right) + \cdots + \frac{1}{m!} \left(1 - \frac{1}{n}\right) \cdots \left(1 - \frac{m-1}{n}\right).$$

Let $n \rightarrow \infty$, keeping m fixed. We get

$$\liminf_{n \rightarrow \infty} t_n \geq 1 + 1 + \frac{1}{2!} + \cdots + \frac{1}{m!},$$

so that

$$s_m \leq \liminf_{n \rightarrow \infty} t_n.$$

Letting $m \rightarrow \infty$, we get

$$e \leq \liminf_{n \rightarrow \infty} t_n.$$

Thus it follows that

$$\limsup_{n \rightarrow \infty} t_n = \liminf_{n \rightarrow \infty} t_n = e,$$

so the desired result follows. □

Lemma 12.39 (Root test). Given $\sum a_n$, let $\alpha = \limsup_{n \rightarrow \infty} \sqrt[n]{|a_n|}$.

- (i) If $\alpha < 1$, $\sum a_n$ converges.
- (ii) If $\alpha > 1$, $\sum a_n$ diverges.
- (iii) If $\alpha = 1$, the test gives no information.

Remark. We use limsup since the limsup of a sequence always exists (in $\overline{\mathbb{R}}$), while the limit may not necessarily exist.

Proof.

- (i) If $\alpha < 1$, choose β such that $\alpha < \beta < 1$. Since $\beta > \limsup_{n \rightarrow \infty} \sqrt[n]{|a_n|}$, there exists $n \in \mathbb{N}$ such that for all

$$n \geq N,$$

$$\sqrt[n]{|a_n|} < \beta,$$

or

$$|a_n| < \beta^n.$$

Note that $\sum \beta^n$ converges since $0 < \beta < 1$. By the comparison test, $\sum a_n$ converges.

(ii) If $\alpha > 1$, $\limsup_{n \rightarrow \infty} \sqrt[n]{|a_n|} > 1$, so there exists a subsequence (a_{n_k}) such that

$$\sqrt[n_k]{|a_{n_k}|} \rightarrow \alpha.$$

Thus $|a_n| > 1$ for infinitely many values of n . Hence $a_n \not\rightarrow 0$, so by the divergence test, $\sum a_n$ diverges.

(iii) Consider the series $\sum \frac{1}{n}$ and $\sum \frac{1}{n^2}$. For each of these series $\alpha = 1$, but the first diverges, the second converges. Hence the condition that $\alpha = 1$ does not give us information on the convergence of a series.

□

Lemma 12.40 (Ratio test). *The series $\sum a_n$*

(i) *converges if $\limsup_{n \rightarrow \infty} \left| \frac{a_{n+1}}{a_n} \right| < 1$;*

(ii) *diverges if $\left| \frac{a_{n+1}}{a_n} \right| \geq 1$ for all $n \geq N_0$ (where N_0 is some fixed integer).*

Proof.

(i) If $\limsup_{n \rightarrow \infty} \left| \frac{a_{n+1}}{a_n} \right| < 1$, there exists $\beta < 1$ and $N \in \mathbb{N}$ such that for all $n \geq N$,

$$\left| \frac{a_{n+1}}{a_n} \right| < \beta.$$

In particular, from $n = N$ to $n = N + k$,

$$\begin{aligned} |a_{N+1}| &< \beta |a_N| \\ |a_{N+2}| &< \beta |a_{N+1}| < \beta^2 |a_N| \\ &\vdots \\ |a_{N+k}| &< \beta^k |a_N| \end{aligned}$$

Hence for all $n \geq N$,

$$\begin{aligned} |a_n| &< |a_N| \beta^{n-N} \\ &= \left(|a_N| \beta^{-N} \right) \beta^n \end{aligned}$$

and taking the sum gives

$$\sum |a_n| < |a_N| \beta^{-N} \sum \beta^n.$$

Since $\beta < 1$, $\sum \beta^n$ converges. By the comparison test, $\sum a_n$ converges.

(ii) Suppose $\left| \frac{a_{n+1}}{a_n} \right| \geq 1$ for all $n \geq N_0$. Then $|a_{n+1}| \geq |a_n|$ for $n \geq N_0$, so $a_n \not\rightarrow 0$. By the divergence test, $\sum a_n$ diverges. □

Remark. The ratio test is easier to apply than the root test (since it is usually easier to compute ratios than n -th roots), but the root test is more powerful, as shown by Theorem 3.37 in [Rud76].

The series $\sum a_n$ is said to *converge absolutely* if the series $\sum |a_n|$ converges.

Lemma 12.41 (Absolute convergence). *If $\sum a_n$ converges absolutely, then $\sum a_n$ converges.*

Proof. Suppose $\sum a_n$ converges absolutely; that is, $\sum |a_n|$ converges. Using the Cauchy criterion, fix $\varepsilon > 0$, there exists $N \in \mathbb{N}$ such that for all $n \geq m \geq N$,

$$\left| \sum_{k=m}^n |a_k| \right| < \varepsilon.$$

Since all the terms are non-negative, we can simply write

$$\sum_{k=m}^n |a_k| < \varepsilon.$$

By the triangle inequality,

$$\left| \sum_{k=m}^n a_k \right| \leq \sum_{k=m}^n |a_k| < \varepsilon.$$

Hence by the Cauchy criterion, $\sum a_n$ converges. □

Note that the converse may not necessarily be true. We say that $\sum a_n$ is *conditionally convergent* if it converges, but does not converge absolutely.

Example 12.42. The alternating harmonic series given by

$$\sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n}$$

converges to $\ln 2$, but it is not absolutely convergent (since the harmonic series diverges).

Summation by Parts

Proposition 12.43 (Partial summation formula). *Given two sequences (a_n) and (b_n) , let the n -partial sum of (a_n) be denoted by*

$$A_n = \sum_{k=0}^n a_k$$

for $n \geq 0$; let $A_{-1} = 0$. Then, if $0 \leq p \leq q$, we have

$$\sum_{n=p}^q a_n b_n = \sum_{n=p}^{q-1} A_n (b_n - b_{n+1}) + A_q b_q - A_{p-1} b_p.$$

Proof. The RHS can be written as

$$\begin{aligned} & \sum_{n=p}^{q-1} A_n b_n + A_q b_q - \sum_{n=p}^{q-1} A_n b_{n+1} - A_{p-1} b_p \\ &= \sum_{n=p}^q A_n b_n - \sum_{n=p-1}^{q-1} A_n b_{n+1} \\ &= \sum_{n=p}^q A_n b_n - \sum_{n=p}^q A_{n-1} b_n \\ &= \sum_{n=p}^q (A_n - A_{n-1}) b_n \\ &= \sum_{n=p}^q a_n b_n \end{aligned}$$

which is equal to the LHS. □

Proposition 12.44. *Suppose (a_n) and (b_n) are sequences such that*

- *the partial sums A_n of $\sum a_n$ form a bounded sequence,*
- *$b_0 \geq b_1 \geq b_2 \geq \dots$,*
- *$b_n \rightarrow 0$.*

Then $\sum a_n b_n = 0$.

Proof. Since the partial sums A_n form a bounded sequence, there exists M such that

$$|A_n| \leq M \quad (\forall n \in \mathbb{N})$$

Since $b_n \rightarrow 0$, fix $\varepsilon > 0$, there exists $N \in \mathbb{N}$ such that

$$b_N \leq \frac{\varepsilon}{2M}.$$

For $q \geq p \geq N$, by the partial summation formula, we have

$$\begin{aligned} \left| \sum_{n=p}^q a_n b_n \right| &= \left| \sum_{n=p}^{q-1} A_n (b_n - b_{n+1}) + A_q b_q - A_{p-1} b_p \right| \\ &\leq M \left| \sum_{n=p}^{q-1} (b_n - b_{n+1}) + b_q + b_p \right| \quad [\because |A_n| \leq M] \\ &= M | (b_p - b_q) + b_q + b_p | = 2M b_p \leq 2M b_n \leq \varepsilon. \end{aligned}$$

By the Cauchy criterion, $\sum a_n b_n$ converges to 0. □

Corollary 12.45 (Alternating series test). *Suppose (c_n) is a sequence such that*

- $|c_1| \geq |c_2| \geq |c_3| \geq \dots$,
- $c_{2m-1} \geq 0, c_{2m} \leq 0$ for $m = 1, 2, 3, \dots$,
- $c_n \rightarrow 0$.

Then $\sum c_n = 0$.

Proof. Let

$$a_n = (-1)^{n+1}, \quad b_n = |c_n|.$$

Note that

- the partial sums of (a_n) are 0s and 1s, so they are bounded;
- $b_0 \geq b_1 \geq b_2 \geq \dots$ holds by assumption;
- $c_n \rightarrow 0$ implies $|c_n| \rightarrow 0$, so $b_n \rightarrow 0$.

Then by 12.44, we have that $\sum a_n b_n = 0$, so $\sum c_n = 0$. □

Addition and Multiplication of Series

Proposition 12.46. If $\sum a_n = A$ and $\sum b_n = B$, then

$$(i) \quad \sum (a_n + b_n) = A + B, \quad (\text{addition})$$

$$(ii) \quad \sum ca_n = cA \text{ for some constant } c. \quad (\text{scalar multiplication})$$

Proof.

(i) Let the n -th partial sums be denoted by

$$A_n = \sum_{k=0}^n a_k, \quad B_n = \sum_{k=0}^n b_k.$$

Then

$$\lim_{n \rightarrow \infty} \sum_{k=0}^n (a_k + b_k) = \lim_{n \rightarrow \infty} (A_n + B_n) = \lim_{n \rightarrow \infty} A_n + \lim_{n \rightarrow \infty} B_n = A + B.$$

(ii) Simply factor out the constant c :

$$\lim_{n \rightarrow \infty} \sum_{k=0}^n ca_k = c \lim_{n \rightarrow \infty} \sum_{k=0}^n a_k = cA.$$

□

The situation becomes more complicated when we consider multiplication of two series. To begin with, we have to define the product. This can be done in several ways; we shall consider the so-called ‘‘Cauchy product’’.

Definition 12.47 (Cauchy product). Given $\sum a_n$ and $\sum b_n$, let

$$c_n = \sum_{k=0}^n a_k b_{n-k} \quad (n = 0, 1, 2, \dots)$$

We call $\sum c_n$ the *product* of the two given series.

This definition may be motivated as follows. If we take two power series $\sum a_n z^n$ and $\sum b_n z^n$, multiply them term by term, and collect terms containing the same power of z , we get

$$\begin{aligned} \left(\sum_{n=0}^{\infty} a_n z^n \right) \left(\sum_{n=0}^{\infty} b_n z^n \right) &= (a_0 + a_1 z + a_2 z^2 + \dots) (b_0 + b_1 z + b_2 z^2 + \dots) \\ &= a_0 b_0 + (a_0 b_1 + a_1 b_0) z + (a_0 b_2 + a_1 b_1 + a_2 b_0) z^2 + \dots \\ &= c_0 + c_1 z + c_2 z^2 + \dots \end{aligned}$$

Setting $z = 1$, we arrive at the above definition.

Note that $\sum c_n$ may not converge, even if $\sum a_n$ and $\sum b_n$ do. However $\sum c_n$ converges if an additional condition is imposed: at least one of the two series converges absolutely.

Proposition 12.48 (Mertens' theorem). *Suppose $\sum a_n = A$, $\sum b_n = B$, and $\sum a_n$ converges absolutely. Then their Cauchy product converges to AB .*

Proof. Let $\sum c_n$ be the Cauchy product of $\sum a_n$ and $\sum b_n$. Let the n -th partial sums be denoted by

$$A_n = \sum_{k=0}^n a_k, \quad B_n = \sum_{k=0}^n b_k, \quad C_n = \sum_{k=0}^n c_k.$$

Also let $\beta_n = B_n - B$. Then

$$\begin{aligned} C_n &= a_0 b_0 + (a_0 b_1 + a_1 b_0) + \cdots + (a_0 b_n + a_1 b_{n-1} + \cdots + a_n b_0) \\ &= a_0 B_n + a_1 B_{n-1} + \cdots + a_n B_0 \\ &= a_0(B + \beta_n) + a_1(B + \beta_{n-1}) + \cdots + a_n(B + \beta_0) \\ &= A_n B + (a_0 \beta_n + a_1 \beta_{n-1} + \cdots + a_n \beta_0) \end{aligned}$$

Our goal is to show that $C_n \rightarrow AB$. Since $A_n B \rightarrow AB$, it suffices to show that

$$\gamma_n = a_0 \beta_n + a_1 \beta_{n-1} + \cdots + a_n \beta_0 \rightarrow 0.$$

We now use the absolute convergence of (a_n) ; let $\alpha = \sum |a_n|$. Fix $\varepsilon > 0$, there exists $N_1 \in \mathbb{N}$ such that

$$n \geq N_1 \implies \sum_{k=0}^n |a_k| - \alpha < \varepsilon$$

since the terms are non-negative. Since $B_n \rightarrow B$, $\beta_n \rightarrow 0$. Then there exists $N_2 \in \mathbb{N}$ such that

$$n \geq N_2 \implies |\beta_n| \leq \varepsilon.$$

Let $N = \max\{N_1, N_2\}$. Then for $n \geq N$, by triangle inequality,

$$\begin{aligned} |\gamma_n| &= |\beta_0 a_n + \cdots + \beta_n a_0| \\ &\leq |\beta_0 a_n + \cdots + \beta_N a_{n-N}| + |\beta_{N+1} a_{n-N-1} + \cdots + \beta_n a_0| \\ &\leq |\beta_0 a_n + \cdots + \beta_N a_{n-N}| + \varepsilon(|a_{n-N-1}| + \cdots + |a_0|) \\ &\leq |\beta_0 a_n + \cdots + \beta_N a_{n-N}| + \varepsilon \alpha. \end{aligned}$$

Keeping N fixed, and letting $n \rightarrow \infty$, we get

$$\limsup_{n \rightarrow \infty} |\gamma_n| \leq \varepsilon \alpha,$$

since $a_n \rightarrow 0$. Since ε is arbitrary, we have $\gamma_n \rightarrow 0$, as desired. □

Proposition 12.49 (Abel's theorem). *Let $\sum a_n = A$, $\sum b_n = B$, $\sum c_n = C$, where $\sum c_n$ is the Cauchy product of $\sum a_n$ and $\sum b_n$. Then $C = AB$.*

Rearrangements

Definition 12.50 (Rearrangement). Let (k_n) be a sequence in which every positive integer appears once and only once. Let

$$a'_n = a_{k_n} \quad (\forall n \in \mathbb{N})$$

We say that $\sum a'_n$ is a *rearrangement* of $\sum a_n$.

If (s_n) and (s'_n) are the sequences of partial sums of (a_n) and (a'_n) respectively, it is easily seen that, in general, these two sequences consist of entirely different numbers. We are thus led to the problem of determining under what conditions all rearrangements of a convergent series will converge and whether the sums are necessarily the same.

Theorem 12.51 (Riemann series theorem). Let $\sum a_n$ be a series of real numbers which converges, but not absolutely. Suppose $-\infty \leq \alpha \leq \beta \leq \infty$. Then there exists a rearrangement $\sum a'_n$ with partial sums s'_n such that

$$\liminf_{n \rightarrow \infty} s'_n = \alpha, \quad \limsup_{n \rightarrow \infty} s'_n = \beta.$$

Proof. Let

$$p_n = \frac{|a_n| + a_n}{2}, \quad q_n = \frac{|a_n| - a_n}{2} \quad (n = 1, 2, \dots).$$

Then $p_n - q_n = a_n$, $p_n + q_n = |a_n|$, $p_n \geq 0$, $q_n \geq 0$.

Claim. The series $\sum p_n$ and $\sum q_n$ must both diverge.

If both were convergent, then

$$\sum (p_n + q_n) = \sum |a_n|$$

would converge, contrary to hypothesis. Since

$$\sum_{n=1}^N a_n = \sum_{n=1}^N (p_n - q_n) = \sum_{n=1}^N p_n - \sum_{n=1}^N q_n,$$

divergence of $\sum p_n$ and convergence of $\sum q_n$ (or vice versa) implies divergence of $\sum a_n$, again contrary to hypothesis.

Now let P_1, P_2, \dots denote the non-negative terms of $\sum a_n$, in the order which they occur, and let Q_1, Q_2, \dots be the absolute values of the negative terms of $\sum a_n$, also in their original order.

The series $\sum P_n$ and $\sum Q_n$ differ from $\sum p_n$ and $\sum q_n$ only by zero terms, and are therefore divergent.

We shall construct sequences (m_n) and (k_n) , such that the series

$$\begin{aligned} & (P_1 + \dots + P_{m_1}) - (Q_1 + \dots + Q_{k_1}) + \\ & (P_{m_1+1} + \dots + P_{m_2}) - (Q_{k_1+1} + \dots + Q_{k_2}) + \dots \end{aligned} \tag{1}$$

which clearly is a rearrangement of $\sum a_n$, satisfies $\liminf_{n \rightarrow \infty} s'_n = \alpha$, $\limsup_{n \rightarrow \infty} s'_n = \beta$.

Choose real-valued sequences (α_n) and (β_n) such that $\alpha_n \rightarrow \alpha$, $\beta_n \rightarrow \beta$, $\alpha_n < \beta_n$, $\beta_1 > 0$.

Let m_1, k_1 be the smallest integers such that

$$\begin{aligned} P_1 + \cdots + P_{m_1} &> \beta_1, \\ P_1 + \cdots + P_{m_1} - (Q_1 + \cdots + Q_{k_1}) &< \alpha_1; \end{aligned}$$

let m_2, k_2 be the smallest integers such that

$$\begin{aligned} (P_1 + \cdots + P_{m_1}) - (Q_1 + \cdots + Q_{k_1}) + (P_{m_1+1} + \cdots + P_{m_2}) &> \beta_2 \\ (P_1 + \cdots + P_{m_1}) - (Q_1 + \cdots + Q_{k_1}) + (P_{m_1+1} + \cdots + P_{m_2}) - (Q_{k_1+1} + \cdots + Q_{k_2}) &< \alpha_2; \end{aligned}$$

and continue in this way. This is possible since $\sum P_n$ and $\sum Q_n$ diverge.

If x_n, y_n denote the partial sums of (1) whose last terms are $P_{m_n}, -Q_{k_n}$, then

$$|x_n - \beta_n| \leq P_{m_n}, \quad |y_n - \alpha_n| \leq Q_{k_n}.$$

Since $P_n \rightarrow 0$ and $Q_n \rightarrow 0$ as $n \rightarrow \infty$, we see that $x_n \rightarrow \beta, y_n \rightarrow \alpha$.

Finally, it is clear that no number less than α or greater than β can be a subsequential limit of the partial sums of (1). to review

□

Theorem 12.52. *If $\sum a_n$ is a series of complex numbers which converges absolutely, then every rearrangement of $\sum a_n$ converges, and they all converge to the same sum.*

Proof. Let $\sum a'_n$ be a rearrangement, with partial sums s'_n . Since $\sum a_n$ converges absolutely, given $\varepsilon > 0$, there exists $N \in \mathbb{N}$ such that $m \geq n \geq N$ implies

$$\sum_{i=n}^m |a_i| < \varepsilon. \tag{1}$$

Now choose p such that the integers $1, 2, \dots, N$ are all contained in the set k_1, \dots, k_p (we use the notation of Definition 12.50). Then if $n > p$, the numbers a_1, \dots, a_N will cancel in the difference $s_n - s'_n$, so that $|s_n - s'_n| < \varepsilon$, by (1). Hence (s'_n) converges to the same sum as (s_n) . to review

□

Exercises

Exercise 12.1. Show the following:

- (i) $\lim_{n \rightarrow \infty} \frac{1}{n^p} = 0$ ($p > 0$)
- (ii) $\lim_{n \rightarrow \infty} \sqrt[p]{p} = 1$ ($p > 0$)
- (iii) $\lim_{n \rightarrow \infty} \sqrt[p]{n} = 1$
- (iv) $\lim_{n \rightarrow \infty} \frac{n^\alpha}{(1+p)^n} = 0$ ($p > 0, \alpha \in \mathbb{R}$)
- (v) $\lim_{n \rightarrow \infty} x^n = 0$ ($|x| < 1$)

Solution.

- (i) Let $\varepsilon > 0$ be given. Take $N = \left\lceil \left(\frac{1}{\varepsilon}\right)^{\frac{1}{p}} \right\rceil + 1$. Then $n \geq N$ implies

$$\left| \frac{1}{n^p} - 0 \right| = \frac{1}{n^p} \leq \frac{1}{N^p} < \frac{1}{\left(\left(\frac{1}{\varepsilon}\right)^{\frac{1}{p}}\right)^p} = \varepsilon.$$

- (ii) We need to consider cases corresponding to different values of p .

Case 1: $p > 1$. Put $x_n = \sqrt[p]{p} - 1$. Then $x_n > 0$, and, by the binomial theorem,

$$1 + nx_n \leq (1 + x_n)^n = p,$$

so that

$$0 < x_n \leq \frac{p-1}{n}.$$

Hence $x_n \rightarrow 0$.

Case 2: $p = 1$. Trivial.

Case 3: $0 < p < 1$. The result is obtained by taking reciprocals.

- (iii) Put $x_n = \sqrt[p]{n} - 1$. Then $x_n \geq 0$, and, by the binomial theorem,

$$n = (1 + x_n)^n \geq \frac{n(n-1)}{2} x_n^2.$$

Hence

$$0 \leq x_n \leq \sqrt{\frac{2}{n-1}} \quad (n \geq 2.)$$

- (iv) Let k be an integer such that $k > \alpha, k > 0$. For $n > 2k$,

$$(1+p)^n > \binom{n}{k} p^k = \frac{n(n-1)\cdots(n-k+1)}{k!} p^k > \frac{n^k p^k}{2^k k!}.$$

Hence

$$0 < \frac{n^\alpha}{(1+p)^n} < \frac{2^k k!}{p^k} n^{\alpha-k} \quad (n > 2k).$$

Since $\alpha - k < 0$, by (i), $n^{\alpha-k} \rightarrow 0$.

(v) Take $\alpha = 0$ in (iv).

□

Exercise 12.2. Let (x_n) be a real sequence, let $\alpha \geq 2$ be a constant. Define the sequence (y_n) as follows:

$$y_n = x_n + \alpha x_{n+1} \quad (n = 1, 2, \dots)$$

Show that if (y_n) is convergent, then (x_n) is also convergent.

Exercise 12.3 ([Rud76] 3.1). Prove that the convergence of (a_n) implies the convergence of $(|a_n|)$. Is the converse true?

Solution. Let $\varepsilon > 0$ be given. Since (a_n) is a Cauchy sequence, there exists $N \in \mathbb{N}$ such that for all $n, m \geq N$,

$$|a_n - a_m| < \varepsilon.$$

See that

$$||a_n| - |a_m|| \leq |a_n - a_m| < \varepsilon,$$

so $(|a_n|)$ is a Cauchy sequence, and therefore must converge.

The converse is not true, as shown by the sequence (a_n) with $a_n = (-1)^n$.

□

Exercise 12.4 ([Rud76] 3.2). Calculate $\lim_{n \rightarrow \infty} (\sqrt{n^2 + n} - n)$.

Solution.

□

Exercise 12.5 ([Rud76] 3.3). The sequence (a_n) is recursively defined by

$$\begin{cases} a_0 = \sqrt{2}, \\ a_{n+1} = \sqrt{2 + a_n} \quad n \geq 0. \end{cases}$$

Show that (a_n) converges.

Solution. We first prove by induction that $a_n \leq a_{n+1} \leq 2$ for all $n \in \mathbb{N}$. For $n = 0$,

$$a_0 = \sqrt{2} \leq \sqrt{2 + \sqrt{2}} = a_1 \leq \sqrt{2 + \sqrt{4}} = 2.$$

If $a_{n-1} \leq a_n \leq 2$, then

$$a_n = \sqrt{2 + a_{n-1}} \leq \sqrt{2 + a_n} = a_{n+1} \leq \sqrt{2 + 2} = 2.$$

Hence (a_n) is monotonically increasing and bounded above by 2. By the monotone convergence theorem, (a_n) converges; let $a_n \rightarrow a$. Applying the limit on both sides of $a_{n+1} = \sqrt{2 + a_n}$,

$$\begin{aligned} \lim_{n \rightarrow \infty} a_{n+1} &= \lim_{n \rightarrow \infty} \sqrt{2 + a_n} \\ a &= \sqrt{2 + a} \\ a &= 2 \text{ or } 1 \end{aligned}$$

Since all $a_n \geq 0$, we must have $a = 2$. □

Exercise 12.6 (Contractive sequence). A complex sequence (x_n) is *contractive* if there exists $k \in [0, 1)$ such that

$$|a_{n+2} - a_{n+1}| \leq k|a_{n+1} - a_n| \quad (\forall n \in \mathbb{N})$$

Show that every contractive sequence is convergent.

Solution. By induction on n , we have

$$|a_{n+1} - a_n| \leq k^{n-1}|a_2 - a_1| \quad (\forall n \in \mathbb{N})$$

Thus

$$\begin{aligned} |a_{n+p} - a_n| &\leq |a_{n+1} - a_n| + |a_{n+2} - a_{n+1}| + \cdots + |a_{n+p} - a_{n+p-1}| \\ &\leq (k^{n-1} + k^n + \cdots + k^{n+p-2})|a_2 - a_1| \\ &\leq k^{n-1} (1 + k + k^2 + \cdots + k^{p-1})|a_2 - a_1| \\ &\leq \frac{k^{n-1}}{1-k}|a_2 - a_1| \end{aligned}$$

for all $n, p \in \mathbb{N}$. Since $k^{n-1} \rightarrow 0$ as $n \rightarrow \infty$ (independently of p), this implies (a_n) is a Cauchy sequence, so it is convergent. □

Exercise 12.7 ([Rud76] 3.4). Find the limit superior and limit inferior of the sequence (a_n) defined by

$$a_1 = 0, \quad a_{2m} = \frac{a_{2m-1}}{2}, \quad a_{2m+1} = a_{2m} + \frac{1}{2}.$$

Solution. We shall prove by induction that

$$a_{2m} = \frac{1}{2} - \frac{1}{2^m}, \quad a_{2m+1} = 1 - \frac{1}{2^m}$$

for $m = 1, 2, \dots$. The second of these equalities is a direct consequence of the first, and so we need only prove the first. Immediate computation shows that $a_2 = 0$ and $a_3 = \frac{1}{2}$. Hence assume that both formulae holds for $m \leq r$. Then

$$a_{2r+2} = \frac{1}{2}a_{2r+1} = \frac{1}{2} \left(1 - \frac{1}{2^r} \right) = \frac{1}{2} - \frac{1}{2^{r+1}}.$$

This completes the induction. We thus have $\limsup_{n \rightarrow \infty} a_n = 1$ and $\liminf_{n \rightarrow \infty} a_n = \frac{1}{2}$. □

Exercise 12.8 ([Rud76] 3.7). Prove that the convergence of $\sum a_n$ implies the convergence of

$$\sum \frac{\sqrt{a_n}}{n}$$

if $a_n \geq 0$.

Exercise 12.9 ([Rud76] 3.8). If $\sum a_n$ converges, and if (b_n) is monotonic and bounded, prove that $\sum a_n b_n$ converges.

Exercise 12.10 ([Rud76] 3.13). Prove that the Cauchy product of two absolutely convergent series converges absolutely.

Exercise 12.11 ([Rud76] 3.23). Suppose (a_n) and (b_n) are Cauchy sequences in a metric space X . Show that the sequence $(d(a_n, b_n))$ converges.

13 Continuity

§13.1 Limit of Functions

Let (X, d_X) and (Y, d_Y) be metric spaces. Let $E \subset X$, then the metric d_X induces a metric on E . Consider a function $f: E \rightarrow Y$. In particular, if $Y = \mathbb{R}$, f is called a *real-valued function*; if $Y = \mathbb{C}$, f is called a *complex-valued function*.

Definition 13.1 (Limit of function). Let p be a limit point of E . We say $\lim_{x \rightarrow p} f(x) = q$ if there exists $q \in Y$ such that

$$\forall \varepsilon > 0, \quad \exists \delta > 0, \quad \forall x \in E, \quad 0 < d_X(x, p) < \delta \implies d_Y(f(x), q) < \varepsilon.$$

The definition conveys the intuitive idea that $f(x)$ can be made arbitrarily close to q by taking x sufficiently close to p .

Remark. Note that $p \in X$, but it is not necessary that $p \in E$ in the above definition. Moreover, even if $p \in E$, we may very well have $f(p) \neq \lim_{x \rightarrow p} f(x)$.

We can recast the above definition in terms of limits of sequences:

Lemma 13.2. Let p be a limit point of E . Then

$$\lim_{x \rightarrow p} f(x) = q \tag{1}$$

if and only if

$$\lim_{n \rightarrow \infty} f(p_n) = q \tag{2}$$

for every sequence (p_n) in $E \setminus \{p\}$ where $p_n \rightarrow p$.

Proof.

\implies Suppose (1) holds. Then fix $\varepsilon > 0$, there exists $\delta > 0$ such that for all $x \in E$,

$$0 < d_X(x, p) < \delta \implies d_Y(f(x), q) < \varepsilon.$$

Let (p_n) be a sequence in $E \setminus \{p\}$. Since $p_n \rightarrow p$, for the same $\varepsilon > 0$, there exists $N \in \mathbb{N}$ such that for all $n \geq N$,

$$0 < d_X(p_n, p) < \delta.$$

This implies that for $n \geq N$, $d_Y(f(p_n), q) < \varepsilon$. Hence by definition $\lim_{n \rightarrow \infty} f(p_n) = q$.

\impliedby Suppose, for a contradiction, (2) holds and (1) does not hold. Then $\lim_{x \rightarrow p} f(x) \neq q$, so

$$\exists \varepsilon > 0, \quad \forall \delta > 0, \quad \exists x \in E, \quad 0 < d_X(x, p) < \delta \quad \text{and} \quad d_Y(f(x), q) \geq \varepsilon.$$

Since (2) holds, taking $\delta_n = \frac{1}{n}$ ($n = 1, 2, \dots$), we thus find a sequence (p_n) in $E \setminus \{p\}$ such that

$$0 < d_X(p_n, p) < \frac{1}{n} \quad \text{and} \quad d_Y(f(p_n), q) \geq \varepsilon.$$

Clearly $p_n \rightarrow p$ but $f(p_n) \not\rightarrow q$, contradicting (2). □

Corollary 13.3. *If f has a limit at p , this limit is unique.*

Proof. Suppose $\lim_{x \rightarrow p} f(x) = q$ and $\lim_{x \rightarrow p} f(x) = q'$. We will show that $q = q'$.

By 13.2, for every sequence (p_n) in $E \setminus \{p\}$ where $p_n \rightarrow p$, we have

$$f(p_n) \rightarrow q \quad \text{and} \quad f(p_n) \rightarrow q'.$$

But the limit of a sequence is unique, so we must have $q = q'$. □

Suppose $f, g: E \rightarrow \mathbb{C}$. Define

$$(f + g)(x) = f(x) + g(x) \quad (x \in E).$$

We define the difference $f - g$, the product fg , and the quotient f/g similarly, with the understanding that the quotient is defined only at $x \in E$ at which $g(x) \neq 0$.

Similarly, if $\mathbf{f}, \mathbf{g}: E \rightarrow \mathbb{R}^k$, we define

$$(\mathbf{f} + \mathbf{g})(x) = \mathbf{f}(x) + \mathbf{g}(x), \quad (\mathbf{f} \cdot \mathbf{g})(x) = \mathbf{f}(x) \cdot \mathbf{g}(x);$$

and if λ is a real number, $(\lambda \mathbf{f})(x) = \lambda \mathbf{f}(x)$.

Lemma 13.4 (Arithmetic properties). *Suppose $E \subset X$, p is a limit point of E . Let $f, g: E \rightarrow \mathbb{C}$, $\lim_{x \rightarrow p} f(x) = A$, $\lim_{x \rightarrow p} g(x) = B$. Then*

$$(i) \quad \lim_{x \rightarrow p} (f + g)(x) = A + B \quad \text{(sum)}$$

$$(ii) \quad \lim_{x \rightarrow p} (fg)(x) = AB \quad \text{(product)}$$

$$(iii) \quad \lim_{x \rightarrow p} \left(\frac{f}{g} \right)(x) = \frac{A}{B} \quad (B \neq 0) \quad \text{(quotient)}$$

Proof. These follow from 13.2 and analogous limit properties of sequences in \mathbb{C} . □

If $\mathbf{f}, \mathbf{g}: E \rightarrow \mathbb{R}^k$, then (i) remains true, and (ii) becomes $\lim_{x \rightarrow p} (\mathbf{f} \cdot \mathbf{g})(x) = \mathbf{A} \cdot \mathbf{B}$.

§13.2 Continuous Functions

Definition 13.5 (Continuity). Suppose $E \subset X$. We say $f: E \rightarrow Y$ is *continuous* at $p \in E$ if

$$\forall \varepsilon > 0, \quad \exists \delta > 0, \quad \forall x \in E, \quad d_X(x, p) < \delta \implies d_Y(f(x), f(p)) < \varepsilon.$$

If f is continuous at every point of E , we say that f is *continuous on E* .

Remark. This definition reflects the intuitive idea that for any arbitrary target distance around $f(p)$, we can always find points $x \in E$ that are sufficiently close to p , such that their images under f are within the target distance around $f(p)$.

Remark. Note that f has to be defined at p in order to be continuous at p . (Compare this with the remark following Definition 13.1.)

Notation. Let X and Y be metric spaces. We denote the space of continuous bounded functions from X to Y as $\mathcal{C}(X, Y)$.

We will often use the next result to show that a function is continuous at a point.

Lemma 13.6. *Let p be a limit point of E . Then f is continuous at p if and only if*

$$\lim_{x \rightarrow p} f(x) = f(p).$$

Proof. Compare Definitions 13.1 and 13.5. □

Corollary 13.7 (Sequential criterion for continuity). *$f: E \subset X \rightarrow Y$ is continuous on E if and only if for every convergent sequence (p_n) in E ,*

$$\lim_{n \rightarrow \infty} f(p_n) = f\left(\lim_{n \rightarrow \infty} p_n\right).$$

Remark. This means that for continuous functions, the limit symbol can be interchanged with the function symbol. Some care is needed in interchanging these symbols because sometimes $(f(p_n))$ converges when (p_n) diverges.

Lemma 13.8. *Let $f, g: X \rightarrow \mathbb{C}$ be continuous on X . Then the following are continuous on X :*

- | | |
|---|------------|
| (i) $f + g$ | (sum) |
| (ii) fg | (product) |
| (iii) $\frac{f}{g}$ ($g(x) \neq 0$ for all $x \in X$) | (quotient) |

Proof. At isolated points of X , there is nothing to prove.

At limit points, the statement follows from 13.4 and 13.6. □

Example 13.9. It is a trivial exercise to show that the following complex-valued functions are continuous on \mathbb{C} :

- constant functions, defined by $f(z) = c$ for all $z \in \mathbb{C}$;
- the identity function, defined by $f(z) = z$ for all $z \in \mathbb{C}$.

Repeated application of the previous result establishes the continuity of every polynomial

$$f(z) = a_0 + a_1z + a_2z^2 + \cdots + a_nz^n$$

where $a_i \in \mathbb{C}$.

We now prove the analogue for Euclidean spaces.

Lemma 13.10.

(i) Let $f_1, \dots, f_k : X \rightarrow \mathbb{R}$, and let $\mathbf{f} : X \rightarrow \mathbb{R}^k$ be defined by

$$\mathbf{f}(x) = (f_1(x), \dots, f_k(x)) \quad (x \in X).$$

Then \mathbf{f} is continuous if and only if each of its components f_1, \dots, f_k is continuous.

(ii) Let $\mathbf{f}, \mathbf{g} : X \rightarrow \mathbb{R}^k$ be continuous on X . Then $\mathbf{f} + \mathbf{g}$ and $\mathbf{f} \cdot \mathbf{g}$ are continuous on X .

Proof. (i) follows from the inequalities

$$|f_j(x) - f_j(y)| \leq |\mathbf{f}(x) - \mathbf{f}(y)| = \left(\sum_{i=1}^k |f_i(x) - f_i(y)|^2 \right)^{1/2}$$

for $j = 1, \dots, k$.

(ii) follows from (i) and 13.8. □

We now consider the composition of functions. The following result shows that a continuous function of a continuous function is continuous.

Proposition 13.11. Suppose X, Y, Z are metric spaces, $E \subset X$. Let

- $f: E \rightarrow Y$,
- $g: f(E) \subset Y \rightarrow Z$,
- $h: E \rightarrow Z$ is defined by $h = g \circ f$.

If f is continuous at $p \in E$, and g is continuous at $f(p)$, then h is continuous at p .

Proof. Let $\varepsilon > 0$ be given. Since g is continuous at $f(p)$, there exists $\eta > 0$ such that for all $y \in f(E)$,

$$d_Y(y, f(p)) < \eta \implies d_Z(g(y), g(f(p))) < \varepsilon. \quad (1)$$

Since f is continuous at p , there exists $\delta > 0$ such that for all $x \in E$,

$$d_X(x, p) < \delta \implies d_Y(f(x), f(p)) < \eta. \quad (2)$$

Combining (1) and (2), it follows that for all $x \in E$,

$$d_X(x, p) < \delta \implies d_Z(h(x), h(p)) = d_Z(g(f(x)), g(f(p))) < \varepsilon.$$

Therefore h is continuous at p . □

Notation. While functions are technically defined on a subset E of a metric space, the complement of E plays no role in the definition of continuity, so we can safely ignore the complement, and think of continuous functions as mappings from one metric space to another.

Continuity and Pre-images of Open or Closed Sets

The following result is another characterisation of continuity.

Lemma 13.12. $f: X \rightarrow Y$ is continuous on X if and only if $f^{-1}(U)$ is open in X for every open set $U \subset Y$.

Proof.

\Rightarrow Suppose f is continuous on X . Let $U \subset Y$ be open. Let $p \in f^{-1}(U)$.

Since $p \in f^{-1}(U)$, there exists $y \in U$ such that $f(p) = y$. By openness of U , there exists $\varepsilon > 0$ such that $B_\varepsilon(y) \subset U$.

Since f is continuous at p , for the same ε , there exists $\delta > 0$ such that for all $x \in X$,

$$d_X(x, p) < \delta \implies d_Y(f(x), y) < \varepsilon,$$

or

$$f(B_\delta(p)) \subset B_\varepsilon(y).$$

Hence

$$B_\delta(p) \subset f^{-1}(f(B_\delta(p))) \subset f^{-1}(B_\varepsilon(y)) \subset f^{-1}(U),$$

so $f^{-1}(U)$ is open in X .

\Leftarrow Suppose $f^{-1}(U)$ is open in X for every open set $U \subset Y$. Fix $p \in X$, let $y = f(p)$. We will show that f is continuous at p .

For every $\varepsilon > 0$, the ball $B_\varepsilon(y)$ is open in Y , so $f^{-1}(B_\varepsilon(y))$ is open in X (by assumption). Now $p \in f^{-1}(B_\varepsilon(y))$, so by openness of $f^{-1}(B_\varepsilon(y))$, there exists $\delta > 0$ such that $B_\delta(p) \subset f^{-1}(B_\varepsilon(y))$. Hence $f(B_\delta(p)) \subset B_\varepsilon(y)$; that is,

$$d_X(x, p) < \delta \implies d_Y(f(x), y) < \varepsilon.$$

Therefore f is continuous at p . □

Corollary 13.13. $f: X \rightarrow Y$ is continuous on X if and only if $f^{-1}(C)$ is closed in X for every closed set $C \subset Y$.

Proof. This follows from the above result, since a set is closed if and only if its complement is open, and since $f^{-1}(E^c) = [f^{-1}(E)]^c$ for every $E \subset Y$. □

Continuity and Compactness

We say $f: E \rightarrow \mathbb{R}^k$ is *bounded* if there exists $M \in \mathbb{R}$ such that $\|f(x)\| \leq M$ for all $x \in E$.

The next result shows that continuous functions preserve compactness.

Proposition 13.14. *Suppose $f: X \rightarrow Y$ is continuous on X , where X is compact. Then $f(X)$ is compact.*

Proof. Let $\{U_i \mid i \in I\}$ be an open cover of $f(X)$. Since f is continuous on X , by 13.12, each of the sets $f^{-1}(U_i)$ is open.

Consider the open cover $\{f^{-1}(U_i) \mid i \in I\}$. Since X is compact, there exist finitely many indices i_1, \dots, i_n such that

$$X \subset \bigcup_{k=1}^n f^{-1}(U_{i_k}).$$

Since $f(f^{-1}(E)) \subset E$ for every $E \subset Y$, we have that

$$f(X) \subset \bigcup_{k=1}^n U_{i_k}.$$

Hence $f(X)$ is compact. □

Corollary 13.15. *If $f: X \rightarrow \mathbb{R}^k$ is continuous on X , where X is compact, then $f(X)$ is closed and bounded. Thus, f is bounded.*

Proof. By 13.14, $f(X)$ is compact. Since $f(X) \subset \mathbb{R}^k$, by the Heine–Borel theorem, $f(X)$ is closed and bounded. □

The result is particularly important when f is a real-valued function; the next result states that a continuous real-valued function on a compact set must attain its minimum and maximum.

Theorem 13.16 (Extreme value theorem). *Suppose $f: X \rightarrow \mathbb{R}$ is continuous, X is compact. Let*

$$M = \sup_{p \in X} f(p), \quad m = \inf_{p \in X} f(p).$$

Then there exist $p, q \in X$ such that $f(p) = M$ and $f(q) = m$.

Proof. From the previous corollary, $f(X)$ is a closed and bounded set in \mathbb{R} . Hence $f(X)$ contains its supremum and infimum, by 11.31. □

Proposition 13.17. *Suppose $f: X \rightarrow Y$ is continuous on X and bijective, X is compact. Then its inverse $f^{-1}: Y \rightarrow X$ is continuous on Y .*

Proof. By 13.12, it suffices to prove that $f(U)$ is open in Y for every open set U in X .

Let U be an open set in X . Then its complement U^c is closed in X . Since U^c is a closed subset of a compact set X , U^c is compact. Thus by 13.14, $f(U^c)$ is a compact subset of Y , so $f(U^c)$ is closed in Y .

Since f is bijective and thus surjective, $f(U)$ is the complement of $f(U^c)$. Hence $f(U)$ is open. \square

Bolzano's Theorem

Lemma 13.18 (Sign-preserving property). *Let $f: [a, b] \rightarrow \mathbb{R}$ be continuous at $c \in [a, b]$, $f(c) \neq 0$. Then there exists $\delta > 0$ such that $f(x)$ has the same sign as $f(c)$ for $c - \delta < x < c + \delta$.*

Proof. Assume $f(c) > 0$. Let $\varepsilon > 0$ be given. By continuity of f , there exists $\delta > 0$ such that

$$c - \delta < x < c + \delta \implies f(c) - \varepsilon < f(x) < f(c) + \varepsilon.$$

Take the δ corresponding to $\varepsilon = \frac{f(c)}{2}$. Then

$$\frac{1}{2}f(c) < f(x) < \frac{3}{2}f(c) \quad (c - \delta < x < c + \delta)$$

so $f(x)$ has the same sign as $f(c)$ for $c - \delta < x < c + \delta$.

The proof is similar if $f(c) < 0$, except that we take $\varepsilon = -\frac{1}{2}f(c)$. □

The next result states that if the graph of $f: [a, b] \rightarrow \mathbb{R}$ lies above the x -axis at a and below the x -axis at b , then the graph must cross the axis somewhere in between. (This should be intuitively obvious.)

Theorem 13.19 (Bolzano). *Suppose $f: [a, b] \rightarrow \mathbb{R}$ is continuous, and $f(a)f(b) < 0$ (that is, $f(a)$ and $f(b)$ have opposite signs). Then there exists $c \in (a, b)$ such that $f(c) = 0$.*

Proof. For definiteness, assume $f(a) > 0$ and $f(b) < 0$. Let

$$A = \{x \in [a, b] \mid f(x) \geq 0\}.$$

Then A is non-empty since $a \in A$, and A is bounded above by b , so A has a supremum in \mathbb{R} ; let $c = \sup A$. Then $a < c < b$.

Claim. $f(c) = 0$.

If $f(c) \neq 0$, by the previous result, there exists $\delta > 0$ such that $f(x)$ has the same sign as $f(c)$ for $c - \delta < x < c + \delta$.

- If $f(c) > 0$, there are points $x > c$ at which $f(x) > 0$, contradicting the definition of c .
- If $f(c) < 0$, then $c - \frac{\delta}{2}$ is an upper bound for A , again contradicting the definition of c .

Therefore we must have $f(c) = 0$. □

Continuity and Connectedness

Proposition 13.20. *Suppose $f: X \rightarrow Y$ is continuous. If $E \subset X$ is connected, then $f(E)$ is connected.*

Proof. We prove the contrapositive. Suppose $f(E)$ is not connected, then $f(E) = A \cup B$ for some $A, B \subset Y$ where $\overline{A} \cap B = \overline{B} \cap A = \emptyset$.

Consider \overline{A} and \overline{B} , which are closed in Y . Since f is continuous, by 13.13, $f^{-1}(\overline{A})$ and $f^{-1}(\overline{B})$ are closed in X ; let $K_A = f^{-1}(\overline{A})$, $K_B = f^{-1}(\overline{B})$. We now want to construct a separation of E .

Let $E_1 = f^{-1}(A) \cap E$, $E_2 = f^{-1}(B) \cap E$. Since $A \cap B = \emptyset$, we have that $E_1 \cap E_2 = \emptyset$. Since $A, B \neq \emptyset$, we have that $E_1, E_2 \neq \emptyset$.

Claim. E_1 and E_2 is a separation of E .

Notice $E_1 \subset K_A$ (which is closed) and $E_2 \subset K_B$ (which is closed). Then $\overline{E_1} \subset K_A$ and $\overline{E_2} \subset K_B$. Note that

$$f^{-1}(\overline{A}) \cap f^{-1}(B) = f^{-1}(\overline{A} \cap B) = \emptyset$$

so $K_A \cap E_2 = \emptyset$. Similarly $K_B \cap E_1 = \emptyset$.

Therefore E is separated. □

The next result says that a continuous real-valued function assumes all intermediate values on an interval.

Theorem 13.21 (Intermediate value theorem). *Suppose $f: [a, b] \rightarrow \mathbb{R}$ is continuous. If $f(a) < f(b)$ and $f(a) < c < f(b)$, then there exists $x \in (a, b)$ such that $f(x) = c$.*

Proof. By 11.62, $[a, b]$ is connected. By the previous result, $f([a, b])$ is a connected subset of \mathbb{R} . Then apply 11.63 and we are done. □

Remark. The converse is not necessarily true. For instance, the *topologist's sine curve*

$$f(x) = \begin{cases} 0 & (x = 0) \\ \sin\left(\frac{1}{x}\right) & (x \neq 0) \end{cases}$$

satisfies the intermediate value property, but f is not continuous.

§13.3 Uniform Continuity

Definition 13.22 (Uniform continuity). We say $f: X \rightarrow Y$ is *uniformly continuous* on X if

$$\forall \varepsilon > 0, \quad \exists \delta > 0, \quad \forall p, q \in X, \quad d_X(p, q) < \delta \implies d_Y(f(p), f(q)) < \varepsilon.$$

Remark. The difference between continuity and uniform continuity is that of one between a local and global property.

- Continuity can be defined at a single point, as δ depends on ε as well as the point p .
- Uniform continuity is a property of a function on a set, as the same δ has to work for *all* $p \in X$ (which ensures a *uniform* rate of closeness across the entire domain.).

Hence uniform continuity is a stronger continuity condition than continuity; a function that is uniformly continuous is continuous but a function that is continuous is not necessarily uniformly continuous.

Example 13.23.

- Let $f(x) = \frac{1}{x}$. Then f is continuous on $(0, 1]$ but not uniformly continuous on $(0, 1]$. To prove this, let $\varepsilon = 10$, and suppose we could find a δ ($0 \leq \delta < 1$) that satisfies the condition of the definition. Taking $p = \delta$, $q = \frac{\delta}{11}$, we obtain $|p - q| < \delta$ and

$$|f(p) - f(q)| = \frac{11}{\delta} - \frac{1}{\delta} = \frac{10}{\delta} > 10.$$

Hence, for these two points we would always have $|f(p) - f(q)| > 10$, contradicting the definition of uniform continuity.

- Let $f(x) = x^2$. Then f is uniformly continuous on $(0, 1]$. To prove this, observe that

$$|f(p) - f(q)| = |p^2 - q^2| = |(p + q)(p - q)| < 2|p - q|.$$

If $|p - q| < \delta$, then $|f(p) - f(q)| < 2\delta$. Hence, for any given ε , we need only take $\delta = \frac{\varepsilon}{2}$ to guarantee that $|f(p) - f(q)| < \varepsilon$ for every $p, q \in (0, 1]$ with $|p - q| < \delta$.

The next result concerns the relationship between continuity and uniform continuity.

Lemma 13.24.

- (i) If $f: X \rightarrow Y$ is uniformly continuous on X , then f is continuous on X .
- (ii) (Heine–Cantor theorem) If $f: X \rightarrow Y$ is continuous on X , and X is compact, then f is uniformly continuous on X .

Proof.

(i)

- (ii) Let $\varepsilon > 0$ be given. Since f is continuous on X , for each $p \in X$, we can associate some $\phi(p) > 0$ such that for all $q \in X$,

$$d_X(p, q) < \phi(p) \implies d_Y(f(p), f(q)) < \frac{\varepsilon}{2}.$$

Consider the collection of open balls centred at each $p \in X$:

$$\left\{ B_{\frac{1}{2}\phi(p)}(p) \mid p \in X \right\}.$$

Since $p \in B_{\frac{1}{2}\phi(p)}(p)$, the above collection of open balls forms an open cover of X . Since X is compact, there exists finitely many points $p_1, \dots, p_n \in X$ such that

$$X \subset \bigcup_{k=1}^n B_{\frac{1}{2}\phi(p_k)}(p_k).$$

Let

$$\delta = \min \left\{ \frac{1}{2}\phi(p_1), \dots, \frac{1}{2}\phi(p_n) \right\}.$$

We claim that this value of δ works in the definition of uniform continuity. Note that $\delta > 0$. (This is one point where the finiteness of the covering, inherent in the definition of compactness, is essential. The minimum of a finite set of positive numbers is positive, whereas the inf of an infinite set of positive numbers may very well be 0.)

Let $p, q \in X$ such that $d_X(p, q) < \delta$. Since X is covered by finitely many open balls, $p \in B_{\frac{1}{2}\phi(p_m)}(p_m)$ for some m ($1 \leq m \leq n$); thus

$$d_X(p, p_m) < \frac{1}{2}\phi(p_m).$$

We also have

$$\begin{aligned} d_X(q, p_m) &\leq d_X(p, q) + d_X(p, p_m) \\ &< \delta + \frac{1}{2}\phi(p_m) \\ &\leq \phi(p_m). \end{aligned}$$

Finally, invoking the continuity of f ,

$$\begin{aligned} d_Y(f(p), f(q)) &\leq d_Y(f(p), f(p_m)) + d_Y(f(q), f(p_m)) \\ &< \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon. \end{aligned}$$

□

Lemma 13.25 (Lebesgue covering lemma). *Suppose $\{U_i \mid i \in I\}$ is an open cover of a compact metric space X . Then there exists $\delta > 0$ such that for all $x \in X$,*

$$B_\delta(x) \subset U_i$$

for some $i \in I$; δ is called a Lebesgue number of the cover.

Proof. Since X is compact, there exist finitely many indices i_1, \dots, i_n such that

$$X \subset \bigcup_{k=1}^n U_{i_k}.$$

For any closed set A , define the distance

$$d(x, A) = \inf_{a \in A} d(x, a).$$

Claim. $d(x, A)$ is a continuous function of x .

Then let the average distance from each x to the complements of U_{i_k} be the function

$$f(x) = \frac{1}{n} \sum_{k=1}^n d(x, U_{i_k}^c).$$

Since f is a sum of continuous functions, f is continuous. Since f is continuous on a compact set, f attains its minimum value; call it δ . See that $\delta > 0$ since $\{U_{i_1}, \dots, U_{i_n}\}$ is an open cover (so $x \in U_{i_k}$ implies $d(x, U_{i_k}^c) > 0$).

For each x , $f(x) \geq \delta$ implies that at least one of the distances $d(x, U_{i_k}^c) \geq \delta$. Hence $B_\delta(x) \subset U_{i_k}$, as desired. \square

§13.4 Discontinuities

We now focus our attention on real-valued functions defined on intervals of the real line.

Definition 13.26 (One-sided limits). Let $f: (a, b) \rightarrow \mathbb{R}$. Let $x \in [a, b)$. The **right-hand limit**, denoted by $f(x+)$ or $\lim_{t \rightarrow x^+} f(t)$, exists if

$$\forall \varepsilon > 0, \quad \exists \delta > 0, \quad x < t < x + \delta < b \implies |f(t) - f(x+)| < \varepsilon.$$

If f is defined at x and if $f(x+) = f(x)$, we say that f is *continuous from the right* at x .

Similarly, let $x \in (a, b]$. The **left-hand limit**, denoted by $f(x-)$ or $\lim_{t \rightarrow x^-} f(t)$, exists if

$$\forall \varepsilon > 0, \quad \exists \delta > 0, \quad a < x - \delta < t < x \implies |f(t) - f(x-)| < \varepsilon.$$

If f is defined at x and if $f(x-) = f(x)$, we say that f is *continuous from the left* at x .

Remark. Compare the above definition with Definition 13.1; for one-sided limits, we are only concerned with half open balls around t (since we only require x to approach t from either the right or left side).

Remark. An equivalent formulation using limits of sequences is presented in [Rud76].

Lemma 13.27. *If $a < x < b$, then f is continuous at c if and only if*

$$f(x) = f(x+) = f(x-).$$

If f is not continuous at x , we say that f is *discontinuous* at x , or that f has a *discontinuity* at x .

Example 13.28 (Dirichlet function). The *Dirichlet function*, defined by

$$f(x) = \begin{cases} 1 & (x \in \mathbb{Q}) \\ 0 & (x \in \mathbb{R} \setminus \mathbb{Q}) \end{cases}$$

is discontinuous everywhere; that is, f is not continuous at any point in \mathbb{R} .

Proof. We consider two cases.

- If $x \in \mathbb{Q}$, then $f(x) = 1$. Take $\varepsilon = \frac{1}{2}$. Since the irrational numbers are dense in the reals, for any $\delta > 0$, we can always find an irrational $y \in \mathbb{R} \setminus \mathbb{Q}$ such that

$$|x - y| < \delta \quad \text{and} \quad |f(x) - f(y)| = 1 \geq \frac{1}{2}.$$

- If $x \in \mathbb{R} \setminus \mathbb{Q}$, then $f(x) = 0$. Again take $\varepsilon = \frac{1}{2}$. Since \mathbb{Q} is dense in \mathbb{R} , for any $\delta > 0$, we can always find $y \in \mathbb{Q}$ such that

$$|x - y| < \delta \quad \text{and} \quad |f(x) - f(y)| = 1 \geq \frac{1}{2}.$$

□

If f is defined on an interval, it is customary to divide discontinuities into two types.

Definition 13.29 (Discontinuities). Let $f: (a, b) \rightarrow \mathbb{R}$. Suppose f is discontinuous at $x \in (a, b)$.

- (i) We say f has a *discontinuity of the first kind* (or a *simple discontinuity*) at x , if $f(x+)$ and $f(x-)$ exist;
- (ii) we say f has a *discontinuity of the second kind* if otherwise.

There are two ways in which a function can have a simple discontinuity: either $f(x+) \neq f(x)$ [in which case the value $f(x)$ is immaterial], or $f(x+) = f(x-) \neq f(x)$.

Example 13.30.

- The function

$$f(x) = \begin{cases} x + 2 & (-3 < x < -2) \\ -x - 2 & (-2 \leq x < 0) \\ x + 2 & (0 \leq x < 1) \end{cases}$$

has a simple discontinuity at $x = 0$, and is continuous at every other point of $(-3, 1)$.

- The Dirichlet function has a discontinuity of the second kind at every $x \in \mathbb{R}$, since both $f(x+)$ and $f(x-)$ do not exist.
- The topologist's sine curve has a discontinuity of the second kind at $x = 0$, since $f(x+)$ does not exist.

§13.5 Monotonic Functions

We now study those functions which never decrease (or never increase) on a given interval.

Definition 13.31 (Monotonicity). $f: (a, b) \rightarrow \mathbb{R}$ is said to be

- (i) *monotonically increasing*, if $f(x_1) \leq f(x_2)$ for any $a < x_1 \leq x_2 < b$;
- (ii) *monotonically decreasing*, if $f(x_1) \geq f(x_2)$ for any $a < x_1 \leq x_2 < b$;
- (iii) **monotonic** if it is either monotonically increasing or monotonically decreasing.

Proposition 13.32. Let $f: (a, b) \rightarrow \mathbb{R}$ be monotonically increasing. Then $f(x+)$ and $f(x-)$ exist for all $x \in (a, b)$; more precisely,

$$\sup_{t \in (a, x)} f(t) = f(x-) \leq f(x) \leq f(x+) = \inf_{t \in (x, b)} f(t).$$

Furthermore, if $a < x < y < b$, then

$$f(x+) \leq f(y-).$$

Analogous results evidently hold for monotonically decreasing functions.

Proof. We will prove the first half of the given statement; the second half can be proven in precisely the same way.

Let $x \in (a, b)$. Since f is monotonically increasing, the set

$$A = \{f(t) \mid a < t < x\}$$

is bounded above by the number $f(x)$. Hence A has a supremum in \mathbb{R} ; let $\alpha = \sup A$. Evidently $\alpha \leq f(x)$.

Claim. $f(x-) = \alpha$.

To prove this, we need to show that for all $\varepsilon > 0$, there exists $\delta > 0$ such that

$$x - \delta < t < x \implies |f(t) - \alpha| < \varepsilon.$$

Let $\varepsilon > 0$ be given. Since $\alpha = \sup A$, there exists $\delta > 0$ such that $a < x - \delta < x$ and

$$\alpha - \varepsilon < f(x - \delta) \leq \alpha. \tag{1}$$

Since f is monotonic, we have

$$f(x - \delta) \leq f(t) \leq \alpha \quad (x - \delta < t < x) \tag{2}$$

Combining (1) and (2) gives

$$|f(t) - \alpha| < \varepsilon \quad (x - \delta < t < x)$$

as desired. Hence $f(x-) = \alpha$.

Next, if $a < x < y < b$, we see from the given statement that

$$f(x+) = \inf_{t \in (x,b)} f(t) = \inf_{t \in (x,y)} f(t)$$

where the last equality is obtained by applying the given statement to (a, y) in place of (a, b) . Similarly,

$$f(y-) = \sup_{t \in (a,y)} f(t) = \sup_{t \in (x,y)} f(t).$$

Comparing these two equations, we conclude that $f(x+) \leq f(y-)$. □

Corollary 13.33. *Monotonic functions have no discontinuities of the second kind.*

Proposition 13.34. *Let $f: (a, b) \rightarrow \mathbb{R}$ be monotonic. Then the set of points of (a, b) at which f is discontinuous is at most countable.*

Proof. Suppose, for the sake of definiteness, that f is monotonically increasing. Let D be the set of points at which f is discontinuous.

For every $x \in D$, we associate a rational number $r(x)$, where

$$f(x-) < r(x) < f(x+).$$

We now check that the rationals picked for two distinct points of discontinuities are different: since $x_1 < x_2$ implies $f(x_1+) \leq f(x_2-)$ (from the previous result), we see that $r(x_1) \neq r(x_2)$ if $x_1 \neq x_2$.

We have thus established a 1-1 correspondence between D and a subset of \mathbb{Q} (which we know is at most countable). Hence D is at most countable. □

§13.6 Lipschitz Continuity

Definition 13.35. $f: X \rightarrow Y$ is *Lipschitz continuous* if there exists $K \geq 0$ such that

$$\forall x, y \in X, \quad d_Y(f(x), f(y)) \leq K d_X(x, y).$$

K is called a *Lipschitz constant* for f ; we also refer to f as *K -Lipschitz*.

Lemma 13.36. *Lipschitz continuity implies uniform continuity.*

Proof. Let $f: X \rightarrow Y$ be K -Lipschitz continuous.

Let $\varepsilon > 0$ be given, let $x, y \in X$. We consider two cases.

Case 1: $K \leq 0$. Then

$$d_X(x, y) \leq 0 d_Y(f(x), f(y))$$

so

$$d_X(x, y) \leq 0 \implies d_X(x, y) = 0 \implies x = y$$

for all $x, y \in X$. Hence f is a constant function, which is uniformly continuous.

Case 2: $K > 0$. Take $\delta = \frac{\varepsilon}{K}$. If $d_X(x, y) < \delta$, then

$$K d_X(x, y) < \varepsilon.$$

By Lipschitz continuity of f ,

$$d_Y(f(x), f(y)) \leq K d_X(x, y).$$

These last two statements together imply $d_Y(f(x), f(y)) < \varepsilon$. Hence f is uniformly continuous on X .

□

We say $f: X \rightarrow Y$ is a *contraction* if it is a K -Lipschitz map for some $K < 1$.

Let $f: X \rightarrow X$, we say $x \in X$ is a *fixed point* if $f(x) = x$.

Theorem 13.37 (Contraction mapping theorem). *Let X be a complete metric space, and $f: X \rightarrow X$ be a contraction. Then f has a unique fixed point.*

Remark. The hypotheses “complete” and “contraction” are necessary. For example, $f: (0, 1) \rightarrow (0, 1)$ defined by $f(x) = Kx$ for any $0 < K < 1$ is a contraction with no fixed point. Also, $f: \mathbb{R} \rightarrow \mathbb{R}$ defined by $f(x) = x + 1$ is not a contraction ($K = 1$) and has no fixed point.

Proof. Pick any $x_0 \in X$. Define a sequence (x_n) by $x_{n+1} = f(x_n)$. Since f is a contraction, we have

$$\begin{aligned} d(x_{n+1}, x_n) &= d(f(x_n), f(x_{n-1})) \\ &\leq Kd(x_n, x_{n-1}) \\ &\leq \dots \\ &\leq K^n d(x_1, x_0) \end{aligned}$$

by induction. Suppose $m \geq n$, then

$$\begin{aligned} d(x_m, x_n) &\leq \sum_{i=n}^{m-1} d(x_{i+1}, x_i) \\ &\leq \sum_{i=n}^{m-1} K^i d(x_1, x_0) \\ &= K^n d(x_1, x_0) \sum_{i=0}^{m-n-1} K^i \\ &\leq K^n d(x_1, x_0) \sum_{i=0}^{\infty} K^i = \frac{K^n}{1-K} d(x_1, x_0). \end{aligned}$$

Thus (x_n) is a Cauchy sequence. Since X is complete, (x_n) converges; let $\lim_{n \rightarrow \infty} x_n = x$ for some $x \in X$.

Claim. x is our unique fixed point.

Note that f is continuous because it is a contraction. Hence

$$f(x) = \lim_{n \rightarrow \infty} f(x_n) = \lim_{n \rightarrow \infty} x_{n+1} = x,$$

so x is a fixed point.

Let x' also be a fixed point. Then

$$d(x, x') = d(f(x), f(x')) = Kd(x, x').$$

As $K < 1$ this means that $d(x, x') = 0$ and hence $x = x'$. The theorem is proved. \square

Note that the proof is constructive. Not only do we know that a unique fixed point exists. We also know how to find it.

§13.7 Infinite Limits and Limits at Infinity

To enable us to operate in the extended real number system, we shall now enlarge the scope of Definition 13.1, reformulating it in terms of open balls.

For any real number x , we have already defined an open ball of x to be any open interval $(x - \delta, x + \delta)$.

Definition 13.38. Let $c \in \mathbb{R}$. A neighbourhood of $+\infty$ is

$$(c, +\infty) := \{x \in \mathbb{R} \mid x > c\}.$$

Similarly, the set $(-\infty, c)$ is a neighbourhood of $-\infty$.

Definition 13.39. Let $f: E \subset \mathbb{R} \rightarrow \mathbb{R}$. We say that $\lim_{t \rightarrow x} f(t) = A$ where A and x are in the extended real number system, if for every neighbourhood U of A there is a neighbourhood V of x such that $V \cap E$ is not empty, and such that $f(t) \in U$ for all $t \in V \cap E, t \neq x$.

Remark. When A and x are real, Definition 13.39 coincides with Definition 13.1.

to do

Lemma 13.40 (Uniqueness of limit). Let $f: E \subset \mathbb{R} \rightarrow \mathbb{R}$. The limit of f at a point x is unique.

Proof. Suppose

$$\lim_{t \rightarrow x} f(t) = A, \quad \lim_{t \rightarrow x} f(t) = A'.$$

We will show that $A' = A$. □

The analogue of Theorem 4.4 is still true, and the proof offers nothing new. We state it, for the sake of completeness.

Lemma 13.41. Let $f, g: E \subset \mathbb{R} \rightarrow \mathbb{R}$. Suppose $\lim_{t \rightarrow x} f(t) = A, \lim_{t \rightarrow x} g(t) = B$. Then

$$(i) \lim_{t \rightarrow x} (f + g)(t) = A + B$$

$$(ii) \lim_{t \rightarrow x} (fg)(t) = AB$$

$$(iii) \lim_{t \rightarrow x} (f/g)(t) = A/B$$

provided the RHS are defined.

Note that $\infty - \infty, 0 \cdot \infty, \infty/\infty, A/0$ are not defined (see Definition 1.23).

Exercises

Exercise 13.1 ([Rud76] 4.1). Suppose $f: \mathbb{R} \rightarrow \mathbb{R}$ satisfies

$$\lim_{h \rightarrow 0} (f(x+h) - f(x-h)) = 0$$

for every $x \in \mathbb{R}$. Does this imply that f is continuous?

Exercise 13.2 ([Rud76] 4.2). If $f: X \rightarrow Y$ is continuous, prove that

$$f(\overline{E}) \subset \overline{f(E)}$$

for every $E \subset X$.

Exercise 13.3 ([Rud76] 4.3). Let $f: X \rightarrow \mathbb{R}$ be continuous. Let the *zero set* of f be

$$Z(f) = \{x \in X \mid f(x) = 0\}.$$

Prove that $Z(f)$ is closed.

Exercise 13.4 ([Rud76] 4.8). Let f be a real uniformly continuous function on the bounded set $E \subset \mathbb{R}$. Prove that f is bounded on E .

Show that the conclusion is false if boundedness of E is omitted from the hypothesis.

Exercise 13.5 ([Rud76] 4.11). Suppose $f: X \rightarrow Y$ is uniformly continuous on X . Prove that $(f(x_n))$ is a Cauchy sequence in Y for every Cauchy sequence (x_n) in X .

Exercise 13.6 ([Rud76] 4.12). A uniformly continuous function of a uniformly continuous function is uniformly continuous.

Exercise 13.7 ([Rud76] 4.14). Let $I = [0, 1]$ be the closed unit interval. Suppose f is a continuous mapping of I into I . Prove that $f(x) = x$ for at least one $x \in I$.

Exercise 13.8 ([Rud76] 4.15). $f: X \rightarrow Y$ is said to be *open* if $f(V)$ is an open set in Y whenever V is an open set in X .

Prove that every continuous open mapping of \mathbb{R} into \mathbb{R} is monotonic.

Exercise 13.9 ([Rud76] 4.16). Let $[x]$ denote the largest integer contained in x , and let $\{x\} = x - [x]$ denote the fractional part of x . What discontinuities do the functions $[x]$ and $\{x\}$ have?

Exercise 13.10 ([Rud76] 4.18). Every rational x can be written in the form $x = \frac{m}{n}$, where $m \in \mathbb{Z}$, $n \in \mathbb{N}$, $\gcd(m, n) = 1$. When $x = 0$, we take $n = 1$. Consider the function f defined on \mathbb{R} by

$$f(x) = \begin{cases} 0 & (x \in \mathbb{R} \setminus \mathbb{Q}) \\ \frac{1}{n} & (x = \frac{m}{n}) \end{cases}$$

Prove that f is continuous at every irrational point, and that f has a simple discontinuity at every rational point.

Exercise 13.11 ([Rud76] 4.26). Suppose X, Y, Z are metric spaces, and Y is compact. Let $f: X \rightarrow Y$, $g: Y \rightarrow Z$ be continuous and injective, and $h = g \circ f$.

Prove that f is uniformly continuous if h is uniformly continuous. *Hint:* g^{-1} has compact domain $g(Y)$, and $f(x) = g^{-1}(h(x))$.

Prove also that f is continuous if h is continuous.

14 Differentiation

§14.1 The Derivative of A Real Function

Definitions and Properties

Definition 14.1 (Derivative). Suppose $f: [a, b] \rightarrow \mathbb{R}$. For any $x \in [a, b]$, if the limit

$$\lim_{t \rightarrow x} \frac{f(t) - f(x)}{t - x} \quad (a < t < b, t \neq x)$$

exists, we call it the *derivative* of f , denoted by f' . If f' is defined at x , we say that f is *differentiable* at x . If f' is defined at every point of $E \subset [a, b]$, we say that f is *differentiable on E* .

f is *continuously differentiable* on E if f' exists at every point of E , and f' is continuous on E .

Lemma 14.2 (Differentiability implies continuity). *If $f: [a, b] \rightarrow \mathbb{R}$ is differentiable at $x \in [a, b]$, then f is continuous at x .*

Proof. Suppose $f: [a, b] \rightarrow \mathbb{R}$ is differentiable at $x \in [a, b]$. Then the limit $\lim_{t \rightarrow x} \frac{f(t) - f(x)}{t - x}$ exists. Thus by arithmetic properties of limits,

$$\begin{aligned} \lim_{t \rightarrow x} [f(t) - f(x)] &= \lim_{t \rightarrow x} \left[\frac{f(t) - f(x)}{t - x} \cdot (t - x) \right] \\ &= \lim_{t \rightarrow x} \frac{f(t) - f(x)}{t - x} \cdot \lim_{t \rightarrow x} (t - x) \\ &= f'(x) \cdot 0 = 0. \end{aligned}$$

Since $\lim_{t \rightarrow x} f(t) = f(x)$, by 13.6, f is continuous at x . □

Remark. The converse is not true; it is easy to construct continuous functions which fail to be differentiable at isolated points.

Example 14.3 (Weierstrass function). Let $0 < a < 1$, let $b > 1$ be an odd integer, and $ab > 1 + \frac{3}{2}\pi$. Then the function

$$W(x) = \sum_{n=0}^{\infty} a^n \cos(b^n \pi x)$$

is continuous and nowhere differentiable on \mathbb{R} .

Example 14.4. One family of pathological examples in calculus is functions of the form

$$f(x) = x^p \sin \frac{1}{x}.$$

For $p = 1$, the function is continuous and differentiable everywhere other than $x = 0$; for $p = 2$, the function is differentiable everywhere, but the derivative is discontinuous.

Lemma 14.5 (Differentiation rules). *Suppose $f, g : [a, b] \rightarrow \mathbb{R}$ are differentiable at $x \in [a, b]$. Then*

(i) *For a constant α , αf is differentiable at x , and* *(scalar multiplication)*

$$(\alpha f)'(x) = \alpha f'(x).$$

(ii) *$f + g$ is differentiable at x , and* *(addition)*

$$(f + g)'(x) = f'(x) + g'(x).$$

(iii) *fg is differentiable at x , and* *(product rule)*

$$(fg)'(x) = f'(x)g(x) + f(x)g'(x).$$

(iv) *f/g (when $g(x) \neq 0$) is differentiable at x , and* *(quotient rule)*

$$\left(\frac{f}{g}\right)'(x) = \frac{f'(x)g(x) - f(x)g'(x)}{g(x)^2}.$$

Proof.

(i)

$$(\alpha f)'(x) = \lim_{t \rightarrow x} \frac{(\alpha f)(t) - (\alpha f)(x)}{t - x} = \alpha \lim_{t \rightarrow x} \frac{f(t) - f(x)}{t - x} = \alpha f'(x).$$

(ii)

$$\begin{aligned} (f \pm g)'(x) &= \lim_{t \rightarrow x} \frac{(f + g)(t) - (f + g)(x)}{t - x} \\ &= \lim_{t \rightarrow x} \frac{f(t) + g(t) - f(x) - g(x)}{t - x} \\ &= \lim_{t \rightarrow x} \frac{f(t) - f(x)}{t - x} + \lim_{t \rightarrow x} \frac{g(t) - g(x)}{t - x} \\ &= f'(x) + g'(x) \end{aligned}$$

(iii)

$$\begin{aligned} (fg)'(x) &= \lim_{t \rightarrow x} \frac{(fg)(t) - (fg)(x)}{t - x} \\ &= \lim_{t \rightarrow x} \frac{f(t)g(t) - f(x)g(x)}{t - x} \\ &= \lim_{t \rightarrow x} \frac{[f(t) - f(x)]g(t) + f(x)[g(t) - g(x)]}{t - x} \\ &= \lim_{t \rightarrow x} \frac{f(t) - f(x)}{t - x} \cdot g(t) + \lim_{t \rightarrow x} f(x) \cdot \frac{g(t) - g(x)}{t - x} \\ &= f'(x)g(x) + f(x)g'(x) \end{aligned}$$

(iv)

$$\begin{aligned}
\left(\frac{f}{g}\right)'(x) &= \lim_{t \rightarrow x} \frac{\left(\frac{f}{g}\right)(t) - \left(\frac{f}{g}\right)(x)}{t - x} \\
&= \lim_{t \rightarrow x} \frac{1}{g(t)g(x)} \left[g(x) \cdot \frac{f(t) - f(x)}{t - x} - f(x) \cdot \frac{g(t) - g(x)}{t - x} \right] \\
&= \frac{f'(x)g(x) - f(x)g'(x)}{g(x)^2}
\end{aligned}$$

□

By induction, we can obtain the following extensions of the differentiation rules.

Corollary. Suppose $f_1, f_2, \dots, f_n : [a, b] \rightarrow \mathbb{R}$ are differentiable at $x \in [a, b]$. Then

(i) $f_1 + f_2 + \dots + f_n$ is differentiable at x , and

$$(f_1 + f_2 + \dots + f_n)'(x) = f_1'(x) + f_2'(x) + \dots + f_n'(x).$$

(ii) $f_1 f_2 \dots f_n$ is differentiable at x , and

$$\begin{aligned}
(f_1 f_2 \dots f_n)'(x) &= f_1'(x) f_2(x) \dots f_n(x) + f_1(x) f_2'(x) \dots f_n(x) \\
&\quad + \dots + f_1(x) f_2(x) \dots f_n'(x).
\end{aligned}$$

The next result concerns the derivative of composition of functions.

Lemma 14.6 (Chain rule). Suppose f is continuous on $[a, b]$, $f'(x)$ exists at $x \in [a, b]$, g is defined on I that contains $f([a, b])$, and g is differentiable at $f(x)$. Then $h = g \circ f$ is differentiable at x , and

$$h'(x) = g'(f(x)) f'(x). \quad (14.1)$$

Proof. By the definition of the derivative, we have

$$f(t) - f(x) = (t - x)[f'(x) + u(t)] \quad (1)$$

$$g(s) - g(f(x)) = (s - f(x))[g'(f(x)) + v(s)] \quad (2)$$

where $t \in [a, b]$, $s \in I$, $\lim_{t \rightarrow x} u(t) = 0$, $\lim_{s \rightarrow f(x)} v(s) = 0$. ($u(t)$ and $v(s)$ can be viewed as some small error terms which eventually go to 0.) Using first (2) and then (1), we obtain

$$\begin{aligned}
h(t) - h(x) &= g(f(t)) - g(f(x)) \\
&= [f(t) - f(x)] \cdot [g'(f(x)) + v(s)] \\
&= (t - x)[f'(x) + u(t)][g'(f(x)) + v(s)],
\end{aligned}$$

or, if $t \neq x$,

$$\frac{h(t) - h(x)}{t - x} = [g'(f(x)) + v(s)][f'(x) + u(t)].$$

Taking limits $t \rightarrow x$, we see that $u(t)$ and $v(s)$ eventually go to 0, so

$$h'(x) = \lim_{t \rightarrow x} \frac{h(t) - h(x)}{t - x} = g'(f(x)) f'(x)$$

as desired. □

Later on when we talk about properties of differentiation such as the intermediate value theorems, we usually have the following requirement on the function:

f is continuous on $[a, b]$, differentiable on (a, b) .

Derivatives of Higher Order

If f has a derivative f' on an interval, and if f' is itself differentiable, we denote the derivative of f' by f'' , and call f'' the *second derivative* of f . Continuing in this manner, we obtain functions

$$f, f', f'', f^{(3)}, f^{(4)}, \dots, f^{(n)},$$

each of which is the derivative of the preceding one. $f^{(n)}$ is called the n -th derivative (or the derivative or order n) of f .

Notation. $C_1[a, b]$ denotes the set of differentiable functions over $[a, b]$ whose derivative is continuous. More generally, $C_n[a, b]$ denotes the set of functions whose n -th derivative is continuous. In particular, $C_0[a, b]$ is the set of continuous functions over $[a, b]$.

§14.2 Mean Value Theorems

Let (X, d) be a metric space.

Definition 14.7 (Local maximum and minimum). $f: X \rightarrow \mathbb{R}$ has

- (i) a **local maximum** at $x_0 \in X$ if there exists $\delta > 0$ such that $f(x_0) \geq f(x)$ for all $x \in B_\delta(x_0)$;
- (ii) a **local minimum** at $x_0 \in X$ if there exists $\delta > 0$ such that $f(x_0) \leq f(x)$ for all $x \in B_\delta(x_0)$.

Our next result is the basis of many applications of differentiation.

Lemma 14.8 (Fermat's theorem). *Suppose $f: [a, b] \rightarrow \mathbb{R}$. If f has a local maximum or minimum at $x \in (a, b)$, and if $f'(x)$ exists, then*

$$f'(x) = 0.$$

Proof. We prove the case for local maxima; the proof for the case for local minima is similar.

Since x is a local maximum, choose $\delta > 0$ such that

$$a < x - \delta < x < x + \delta < b,$$

and $f(x) \geq f(t)$ for all $x - \delta < t < x + \delta$.

- If $x - \delta < t < x$, then

$$\frac{f(t) - f(x)}{t - x} \geq 0.$$

Letting $t \rightarrow x$, we see that $f'(x) \geq 0$.

- If $x < t < x + \delta$, then

$$\frac{f(t) - f(x)}{t - x} \leq 0.$$

Letting $t \rightarrow x$, we see that $f'(x) \leq 0$.

Hence $f'(x) \geq 0$. □

Theorem 14.9 (Rolle's theorem). *Suppose f is continuous on $[a, b]$ and differentiable in (a, b) . If $f(a) = f(b)$, then there exists $c \in (a, b)$ such that*

$$f'(c) = 0.$$

The idea is to show that f has a local maximum/minimum, then by Fermat's theorem this will then be the stationary point that we're trying to find.

Proof. Since f is continuous on $[a, b]$, by the extreme value theorem (13.16), f attains its maximum M and minimum m .

- If M and m both equal $f(a) = f(b)$, then f is simply a constant function; hence $f'(x) = 0$ for all $x \in [a, b]$.

- Otherwise, f has a maximum/minimum that does not equal $f(a) = f(b)$. Then there exists $c \in (a, b)$ such that $f(c)$ is a local maximum/minimum. Since f is differentiable on (a, b) , $f'(c)$ exists, so by Fermat's theorem, $f'(c) = 0$.

□

Theorem 14.10 (Generalised mean value theorem). *Suppose f and g are continuous on $[a, b]$ and differentiable in (a, b) . Then there exists $c \in (a, b)$ such that*

$$[f(b) - f(a)]g'(c) = [g(b) - g(a)]f'(c). \quad (14.2)$$

Proof. For $t \in [a, b]$, consider the *auxilliary function*

$$h(t) = [f(b) - f(a)]g(t) - [g(b) - g(a)]f(t).$$

Then h is continuous on $[a, b]$, and h is differentiable on (a, b) . Moreover,

$$h(a) = f(b)g(a) - f(a)g(b) = h(b).$$

By Rolle's theorem, there exists $c \in (a, b)$ such that $h'(c) = 0$; that is,

$$[f(b) - f(a)]g'(c) = [g(b) - g(a)]f'(c)$$

as desired. □

Theorem 14.11 (Mean value theorem). *Suppose f is continuous on $[a, b]$ and differentiable in (a, b) . Then there exists $c \in (a, b)$ such that*

$$f(b) - f(a) = f'(c)(b - a). \quad (14.3)$$

Proof. Take $g(x) = x$ in 14.10. □

Lemma 14.12. *Suppose f is differentiable in (a, b) .*

- (i) *If $f'(x) \geq 0$ for all $x \in (a, b)$, then f is monotonically increasing.*
- (ii) *If $f'(x) = 0$ for all $x \in (a, b)$, then f is constant.*
- (iii) *If $f'(x) \leq 0$ for all $x \in (a, b)$, then f is monotonically decreasing.*

Proof. All conclusions can be read off from the equation

$$f'(x) = \frac{f(x_2) - f(x_1)}{x_2 - x_1},$$

which is valid, for each pair of numbers x_1, x_2 in (a, b) , for some x between x_1 and x_2 . □

§14.3 Continuity of Derivatives

The following result implies some sort of a “intermediate value” property of derivatives that is similar to continuous functions.

Theorem 14.13 (Darboux’s theorem). *Suppose f is differentiable on $[a, b]$, and suppose $f'(a) < c < f'(b)$. Then there exists $x \in (a, b)$ such that $f'(x) = c$.*

Proof. For $t \in (a, b)$, consider the auxilliary function

$$g(t) = f(t) - ct.$$

Then

$$g'(a) = f'(a) - c < 0,$$

so there exists $t_1 \in (a, b)$ such that $g(t_1) < g(a)$. Similarly,

$$g'(b) = f'(b) - c > 0,$$

so there exists $t_2 \in (a, b)$ such that $g(t_2) < g(b)$.

By the extreme value theorem, g attains its minimum on $[a, b]$. From above, $g(a)$ and $g(b)$ cannot be minimums, so g attains its minimum at $x \in (a, b)$. By Fermat’s theorem, $g'(x) = 0$. Hence $f'(x) = c$, as desired. \square

Corollary 14.14. *If f is differentiable on $[a, b]$, then f' cannot have any simple discontinuities on $[a, b]$.*

§14.4 L'Hopital's Rule

The following result is frequently used in the evaluation of limits.

Lemma 14.15 (L'Hopital's rule). *Suppose f and g are differentiable over (a, b) , with $g'(x) \neq 0$ for all $x \in (a, b)$, where $-\infty \leq a < b \leq +\infty$. If either*

$$(i) \lim_{x \rightarrow a} f(x) = 0 \text{ and } \lim_{x \rightarrow a} g(x) = 0; \text{ or}$$

$$(ii) \lim_{x \rightarrow a} g(x) = +\infty,$$

and

$$\lim_{x \rightarrow a} \frac{f'(x)}{g'(x)} = A,$$

then

$$\lim_{x \rightarrow a} \frac{f(x)}{g(x)} = A.$$

The analogous statement is of course also true if $x > b$, or if $g(x) \rightarrow -\infty$ in (ii).

Note that we now use the limit concept in the extended sense of Definition 13.39.

Proof. We first consider the case in which $-\infty \leq A < +\infty$. Choose $q \in \mathbb{R}$ such that $A < q$, and choose $r \in \mathbb{R}$ such that $A < r < q$. By (13) there exists $c \in (a, b)$ such that $a < x < c$ implies

$$\frac{f'(x)}{g'(x)} < r.$$

If $a < x < y < c$, then by the generalised mean value theorem (14.10), there exists $t \in (x, y)$ such that

$$\frac{f(x) - f(y)}{g(x) - g(y)} = \frac{f'(t)}{g'(t)} < r.$$

(i) Suppose $\lim_{x \rightarrow a} f(x) = 0$ and $\lim_{x \rightarrow a} g(x) = 0$. Let $x \rightarrow a$ in (18), we see that

$$\frac{f(y)}{g(y)} \leq r < q \quad (a < y < c).$$

(ii) Next, suppose $\lim_{x \rightarrow a} g(x) = +\infty$. Keeping y fixed in (18), we can choose a point $c_1 \in (a, y)$ such that $g(x) > g(y)$ and $g(x) > 0$ if $a < x < c_1$. Multiplying (18) by $[g(x) - g(y)]/g(x)$, we obtain

$$\frac{f(x)}{g(x)} < r - r \frac{g(y)}{g(x)} + \frac{f(y)}{g(x)} \quad (a < x < c_1).$$

If we let $x \rightarrow a$ in (20), (15) shows that there exists $c_2 \in (a, c_1)$ such that

$$\frac{f(x)}{g(x)} < q \quad (a < x < c_2).$$

Summing up, (19) and (21) show that for any q , subject only to the condition $A < q$, there is a point c_2 such that $f(x)/g(x) < q$ if $a < x < c_2$.

In the same manner, if $-\infty < A \leq +\infty$, and p is chosen so that $p < A$, we can find a point c_3 such that

$$p < \frac{f(x)}{g(x)} \quad (a < x < c_3),$$

and (16) follows from these two statements. _____

□

to review
proof

§14.5 Taylor's Theorem

Theorem 14.16 (Taylor's theorem). *Suppose $f: [a, b] \rightarrow \mathbb{R}$, $f^{(n-1)}$ is continuous on $[a, b]$, $f^{(n)}$ exists on (a, b) . Assume that $c \in [a, b]$. Let the Taylor polynomial of degree $n - 1$ of f at $x = c$ be*

$$\begin{aligned} P_{n-1}(x) &= \sum_{k=0}^{n-1} \frac{f^{(k)}(c)}{k!} (x - c)^k \\ &= f(c) + f'(c)(x - c) + \frac{f''(c)}{2!} (x - c)^2 + \cdots + \frac{f^{(n-1)}(c)}{(n-1)!} (x - c)^{n-1}. \end{aligned}$$

Then for every $x \in [a, b]$, $x \neq c$, there exists z_x between x and c such that

$$f(x) = P_{n-1}(x) + \frac{f^{(n)}(z_x)}{n!} (x - c)^n. \quad (14.4)$$

For $n = 1$, this is just the mean value theorem. In general, the theorem shows that f can be approximated by a polynomial of degree $n - 1$, and that Eq. (14.4) allows us to accurately estimate the error.

Proof. Let M be the number defined by

$$f(x) = P_{n-1}(x) + M(x - c)^n.$$

We claim that $n!M = f^{(n)}(z_x)$ for some z_x between x and c .

For all $x \in [a, b]$, let

$$g(x) = f(x) - P_{n-1}(x) - M(x - c)^n.$$

Then for all $x \in (a, b)$,

$$g^{(n)}(x) = f^{(n)}(x) - n!M.$$

Hence our proof will be complete if we can show that $g^{(n)}(z_x) = 0$ for some z_x between c and x .

Since $P_{n-1}^{(k)}(c) = f^{(k)}(c)$ for $k = 0, \dots, n - 1$, we have

$$g(c) = g'(c) = \cdots = g^{(n-1)}(c) = 0.$$

By our choice of M , we have that $g(x) = 0$. By the mean value theorem, there exists x_1 between x and c such that $g'(x_1) = 0$. Since $g'(c) = 0$, we conclude similarly that $g''(x_2) = 0$ for some x_2 between x_1 and c . After n steps we arrive at the conclusion that $g^{(n)}(x_n) = 0$ for some x_n between x_{n-1} and c , that is, between x and c . \square

Example 14.17.

$$\begin{aligned}e^x &= 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots \\ \sin x &= x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots \\ \cos x &= 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots \\ \ln(1+x) &= x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots\end{aligned}$$

There's a lot of things to say about these equations, for example the one for $\ln(1+x)$ only works for $|x| < 1$. Also, if we want the RHS of the expression to be an infinite power series, $f(x)$ has to be smooth (infinitely differentiable).

Exercises

Exercise 14.1. Let f and g be continuous on $[a, b]$ and differentiable on (a, b) . If $f'(x) = g'(x)$, then $f(x) = g(x) + C$.

Exercise 14.2. Given that $f(x) = x^\alpha$ where $0 < \alpha < 1$. Prove that f is uniformly continuous on $[0, +\infty)$.

Exercise 14.3. Let f be continuous on $[0, 1]$ and differentiable on $(0, 1)$ where $f(0) = f(1) = 0$. Prove that there exists $c \in (0, 1)$ such that

$$f(x) + f'(x) = 0.$$

15 Riemann–Stieltjes Integral

The present chapter is based on a definition of the Riemann integral which depends very explicitly on the order structure of the real line. Accordingly, we begin by discussing integration of real-valued functions on intervals. Extensions to complex- and vector-valued functions on intervals follow in later sections.

§15.1 Definition of Riemann–Stieltjes Integral

To approximate the area under the curve of a function, we partition the interval into finitely many sub-intervals, then multiply the width of each sub-interval by its height.

- For the height, we can choose to either use the supremum of the function over the interval or the infimum. Obviously, using the supremum will provide an upper bound, and using the infimum will provide a lower bound.
- For the width, we use the difference between the two endpoints in their output values when input into a monotonically increasing function α .

The upper Riemann integral is the infimum of upper bounds over all possible partitions. The lower Riemann integral is similarly defined. If they are equal, then the function is said to be Riemann–Stieltjes integrable.

Notation and Preliminaries

A **partition** P of a closed interval $[a, b]$ is a finite set of points $\{x_0, x_1, \dots, x_n\}$, where

$$a = x_0 \leq x_1 \leq \dots \leq x_{n-1} \leq x_n = b.$$

Notation. Denote the set of all partitions of $[a, b]$ by $\mathcal{P}[a, b]$.

Let $f: [a, b] \rightarrow \mathbb{R}$ be bounded. Denote

$$M_i = \sup_{x \in [x_{i-1}, x_i]} f(x), \quad m_i = \inf_{x \in [x_{i-1}, x_i]} f(x) \quad (i = 1, \dots, n).$$

Let α be a monotonically increasing function on $[a, b]$. Denote

$$\Delta\alpha_i = \alpha(x_i) - \alpha(x_{i-1}) \quad (i = 1, \dots, n).$$

(These suprema and infima are well-defined, finite real numbers due to the boundedness of f .)

The *upper sum* and *lower sum* of f with respect to the partition P and α are respectively

$$U(f, \alpha; P) = \sum_{i=1}^n M_i \Delta \alpha_i,$$

$$L(f, \alpha; P) = \sum_{i=1}^n m_i \Delta \alpha_i.$$

insert
diagram

The partition P' is a **refinement** of P if $P' \supset P$. Given two partitions P_1 and P_2 , we say that P' is their *common refinement* if $P' = P_1 \cup P_2$.

Intuitively, a refinement will give a better estimation than the original partition, so the upper and lower sums of a refinement should be more restrictive.

Lemma 15.1. *If P' is a refinement of P , then*

(i) $L(f, \alpha; P) \leq L(f, \alpha; P')$

(ii) $U(f, \alpha; P') \leq U(f, \alpha; P)$

Proof.

- (i) Suppose first that P' contains just one point more than P . Let this extra point be x' , and suppose $x_{i-1} < x' < x_i$ for some i ($1 \leq i \leq n$), where $x_{i-1}, x_i \in P$. Let

$$w_1 = \inf_{x \in [x_{i-1}, x']} f(x), \quad w_2 = \inf_{x \in [x', x_i]} f(x).$$

Let, as before,

$$m_i = \inf_{x \in [x_{i-1}, x_i]} f(x).$$

Clearly $w_1 \geq m_i$ and $w_2 \geq m_i$. Then

$$\begin{aligned} &L(f, \alpha; P') - L(f, \alpha; P) \\ &= w_1 (\alpha(x') - \alpha(x_{i-1})) + w_2 (\alpha(x_i) - \alpha(x')) - m_i (\alpha(x_i) - \alpha(x_{i-1})) \\ &= \underbrace{(w_1 - m_i)}_{\geq 0} \underbrace{(\alpha(x') - \alpha(x_{i-1}))}_{> 0} + \underbrace{(w_2 - m_i)}_{\geq 0} \underbrace{(\alpha(x_i) - \alpha(x'))}_{> 0} \\ &\geq 0 \end{aligned}$$

and hence $L(f, \alpha; P) \leq L(f, \alpha; P')$.

If P' contains k more points than P , we repeat this reasoning k times.

- (ii) Analogous to the proof of (i).

□

Since f is bounded, there exist m and M such that $m \leq f(x) \leq M$ for all $x \in [a, b]$. Hence for every partition P ,

$$m (\alpha(b) - \alpha(a)) \leq L(f, \alpha; P) \leq U(f, \alpha; P) \leq M (\alpha(b) - \alpha(a))$$

so that the numbers $L(f, \alpha; P)$ and $U(f, \alpha; P)$ form a bounded set. This shows that the upper and lower integrals are defined for every bounded function f . We now define the *upper and lower Riemann–Stieltjes*

integrals respectively as

$$\int_a^{\bar{b}} f \, d\alpha := \inf_{P \in \mathcal{P}[a,b]} U(f, \alpha; P)$$

$$\int_a^{\underline{b}} f \, d\alpha := \sup_{P \in \mathcal{P}[a,b]} L(f, \alpha; P)$$

where we take inf and sup over all partitions.

One would expect the lower RS integral to be less than or equal to the upper RS integral. We now show this.

Lemma 15.2.

$$\int_a^{\underline{b}} f \, d\alpha \leq \int_a^{\bar{b}} f \, d\alpha .$$

Proof. Let P' be the common refinement of partitions P_1 and P_2 ; that is, $P' = P_1 \cup P_2$. Clearly $P' \supset P_1$; by 15.1,

$$L(f, \alpha; P_1) \leq L(f, \alpha; P').$$

Similarly, $P' \supset P_2$, so

$$U(f, \alpha; P') \leq U(f, \alpha; P_2).$$

Clearly $L(f, \alpha; P') \leq U(f, \alpha; P')$. Thus combining the above two equations gives

$$L(f, \alpha; P_1) \leq U(f, \alpha; P_2).$$

Fix P_2 and take sup over all P_1 gives

$$\int_a^{\underline{b}} f \, d\alpha \leq U(f, \alpha; P_2).$$

Then taking inf over all P_2 gives

$$\int_a^{\underline{b}} f \, d\alpha \leq \int_a^{\bar{b}} f \, d\alpha .$$

□

Defining the Integral

Definition 15.3 (Riemann–Stieltjes integral). We say $f : [a, b] \rightarrow \mathbb{R}$ is *Riemann–Stieltjes integrable* with respect to α over $[a, b]$, if

$$\int_a^b f \, d\alpha = \int_a^{\bar{b}} f \, d\alpha.$$

We call the common value the *Riemann–Stieltjes integral* of f with respect to α over $[a, b]$, and denote it as

$$\int_a^b f \, d\alpha.$$

The functions f and α are referred to as the *integrand* and the *integrator*, respectively.

Notation. $\mathcal{R}(\alpha)$ denotes the set of Riemann–Stieltjes integrable functions with respect to α .

In particular, when $\alpha(x) = x$, we call the corresponding Riemann–Stieltjes integral the *Riemann integral*, and use \mathcal{R} to denote the set of Riemann integrable functions.

Notation. Since x is a “dummy variable” and may be replaced by any other variable, we shall omit it.

Example 15.4 (Dirichlet function). The *Dirichlet function* is defined over $[0, 1]$ by

$$f(x) = \begin{cases} 1 & (x \in \mathbb{Q}) \\ 0 & (x \in \mathbb{R} \setminus \mathbb{Q}) \end{cases}$$

For each subinterval $[x_{i-1}, x_i]$, due to the density of rationals and irrationals, $[x_{i-1}, x_i]$ contains both rationals and irrationals, so $M_i = 1$ and $m_i = 0$. Thus for any partition P ,

$$U(f; P) = 1, \quad L(f; P) = 0.$$

Therefore,

$$1 = \int_a^{\bar{b}} f \, d\alpha \neq \int_a^b f \, d\alpha = 0$$

so the Dirichlet function is not Riemann–Stieltjes integrable.

The next result is particularly useful in determining the Riemann–Stieltjes integrability of a function. We will use it many times later.

Lemma 15.5 (Integrability criterion). $f \in \mathcal{R}(\alpha)$ if and only if

$$\forall \varepsilon > 0, \quad \exists P, \quad U(f, \alpha; P) - L(f, \alpha; P) < \varepsilon.$$

Proof.

\Rightarrow Suppose $f \in \mathcal{R}(\alpha)$. Let $\varepsilon > 0$ be given. Then there exists partitions P_1 and P_2 such that

$$U(f, \alpha; P_2) - \int_a^b f \, d\alpha < \frac{\varepsilon}{2}$$

and

$$\int_a^b f \, d\alpha - L(f, \alpha; P_1) < \frac{\varepsilon}{2}.$$

Choose P to be the common refinement of P_1 and P_2 . Then

$$\begin{aligned} U(f, \alpha; P) &\leq U(f, \alpha; P_2) \\ &< \int_a^b f \, d\alpha + \frac{\varepsilon}{2} \\ &< L(f, \alpha; P_1) + \varepsilon \\ &\leq L(f, \alpha; P) + \varepsilon. \end{aligned}$$

Hence for this partition P , we have

$$U(f, \alpha; P) - L(f, \alpha; P) < \varepsilon.$$

$\square \Leftarrow$ By 15.2, for every partition P ,

$$L(f, \alpha; P) \leq \int_a^b f \, d\alpha \leq \int_a^{\bar{b}} f \, d\alpha \leq U(f, \alpha; P).$$

Since $U(f, \alpha; P) - L(f, \alpha; P) < \varepsilon$, we have

$$0 \leq \int_a^{\bar{b}} d\alpha - \int_a^b f \, d\alpha < \varepsilon.$$

Since this holds for all $\varepsilon > 0$, we have

$$\int_a^{\bar{b}} f \, d\alpha = \int_a^b f \, d\alpha.$$

Hence $f \in \mathcal{R}(\alpha)$. \square

Useful Identities

Proposition 15.6 (Cauchy criterion).

(i) If $U(f, \alpha; P) - L(f, \alpha; P) < \varepsilon$ holds for some P and some $\varepsilon > 0$, then $U(f, \alpha; P') - L(f, \alpha; P') < \varepsilon$ holds (with the same ε) for every refinement of P , P' .

(ii) If $U(f, \alpha; P) - L(f, \alpha; P) < \varepsilon$ holds for $P = \{x_0, \dots, x_n\}$, and

$$s_i, t_i \in [x_{i-1}, x_i] \quad (i = 1, \dots, n)$$

then

$$\sum_{i=1}^n |f(s_i) - f(t_i)| \Delta\alpha_i < \varepsilon.$$

(iii) If $f \in \mathcal{R}(\alpha)$ and the hypotheses of (ii) hold, then

$$\left| \sum_{i=1}^n f(t_i) \Delta\alpha_i - \int_a^b f \, d\alpha \right| < \varepsilon.$$

Proof.

(i) Suppose $U(f, \alpha; P) - L(f, \alpha; P) < \varepsilon$ holds for some partition P and some $\varepsilon > 0$. By 15.1, for any refinement P' ,

$$U(f, \alpha; P') \leq U(f, \alpha; P), \quad L(f, \alpha; P') \leq L(f, \alpha; P).$$

Hence

$$U(f, \alpha; P') - L(f, \alpha; P') \leq U(f, \alpha; P) - L(f, \alpha; P) < \varepsilon.$$

(ii) Since

$$f(s_i), f(t_i) \in [m_i, M_i] \quad (i = 1, \dots, n)$$

it follows that

$$|f(s_i) - f(t_i)| \leq M_i - m_i.$$

Hence

$$\sum_{i=1}^n |f(s_i) - f(t_i)| \Delta\alpha_i \leq U(f, \alpha; P) - L(f, \alpha; P) < \varepsilon.$$

(iii) The desired result follows from the two inequalities

$$L(f, \alpha; P) \leq \sum_{i=1}^n f(t_i) \Delta\alpha_i \leq U(f, \alpha; P)$$

$$L(f, \alpha; P) \leq \int_a^b f \, d\alpha \leq U(f, \alpha; P)$$

□

The next result states that all continuous functions are integrable.

Proposition 15.7 (Continuity implies integrability). *If f is continuous on $[a, b]$, then $f \in \mathcal{R}(\alpha)$.*

Proof. Let $\varepsilon > 0$ be given. Choose $\eta > 0$ such that

$$(\alpha(b) - \alpha(a))\eta < \varepsilon.$$

Since f is continuous on $[a, b]$ which is compact, by 13.24, f is uniformly continuous on $[a, b]$. Thus there exists $\delta > 0$ such that for all $x, y \in [a, b]$,

$$|x - y| < \delta \implies |f(x) - f(y)| < \eta.$$

If P is any partition of $[a, b]$ such that $\Delta x_i < \delta$ for $i = 1, \dots, n$, then

$$M_i - m_i \leq \eta \quad (i = 1, \dots, n).$$

Hence

$$\begin{aligned} U(f, \alpha; P) - L(f, \alpha; P) &= \sum_{i=1}^n (M_i - m_i) \Delta \alpha_i \\ &\leq \eta \sum_{i=1}^n \Delta \alpha_i = \eta (\alpha(b) - \alpha(a)) < \varepsilon. \end{aligned}$$

Therefore $f \in \mathcal{R}(\alpha)$, by the integrability criterion (15.5). □

Proposition 15.8. *If f is monotonic on $[a, b]$, and if α is continuous on $[a, b]$, then $f \in \mathcal{R}(\alpha)$.*

Proof. Let $\varepsilon > 0$ be given. For any positive integer n , choose a partition P such that

$$\Delta \alpha_i = \frac{\alpha(b) - \alpha(a)}{n} \quad (i = 1, \dots, n).$$

This is possible by the intermediate value theorem, due to the continuity of α .

Suppose that f is monotonically increasing (the proof is analogous in the other case). Then

$$M_i = f(x_i), \quad m_i = f(x_{i-1}) \quad (i = 1, \dots, n).$$

Hence

$$\begin{aligned} U(f, \alpha; P) - L(f, \alpha; P) &= \sum_{i=1}^n (M_i - m_i) \Delta \alpha_i \\ &= \frac{\alpha(b) - \alpha(a)}{n} \sum_{i=1}^n (f(x_i) - f(x_{i-1})) \\ &= \frac{\alpha(b) - \alpha(a)}{n} (f(b) - f(a)) < \varepsilon \end{aligned}$$

if n is taken large enough. Hence $f \in \mathcal{R}(\alpha)$, by the integrability criterion. □

Proposition 15.9. *Suppose f is bounded on $[a, b]$, f has only finitely many points of discontinuity on $[a, b]$, and α is continuous at every point at which f is discontinuous. Then $f \in \mathcal{R}(\alpha)$.*

Proof. Let $\varepsilon > 0$ be given. Since f is bounded, let $M = \sup |f(x)|$. Let E be the set of points at which f is discontinuous.

Since E is finite, and α is continuous at every point of E , we can cover E by finitely many disjoint intervals $[u_j, v_j] \subset [a, b]$ such that the sum of the corresponding differences $\sum_j (\alpha(v_j) - \alpha(u_j)) < \varepsilon$. Furthermore, we can place these intervals in such a way that every point of $E \cap (a, b)$ lies in the interior of some $[u_j, v_j]$.

Remove the segments (u_j, v_j) from $[a, b]$. The remaining set K is compact. Hence f is uniformly continuous on K , so there exists $\delta > 0$ such that for all $s, t \in K$,

$$|s - t| < \delta \implies |f(s) - f(t)| < \varepsilon.$$

Now form a partition $P = \{x_0, x_1, \dots, x_n\}$ of $[a, b]$ as follows: Each u_j occurs in P . Each v_j occurs in P . No point of any segment (u_j, v_j) occurs in P . If x_{i-1} is not one of the u_j , then $\Delta x_i < \delta$.

Note that $M_i - m_i \leq 2M$ for every i , and that $M_i - m_i < \varepsilon$ unless x_{i-1} is one of the u_j . Hence

$$\begin{aligned} U(f, \alpha; P) - L(f, \alpha; P) &= \sum_{i=1}^n (M_i - m_i) \Delta \alpha_i \\ &\leq (\alpha(b) - \alpha(a)) \varepsilon + 2M\varepsilon. \end{aligned}$$

Since ε is arbitrary, we have $f \in \mathcal{R}(\alpha)$, by the integrability criterion. □

The next result states that a uniformly continuous function of an integrable function is also integrable.

Proposition 15.10. *Suppose $f \in \mathcal{R}(\alpha)$, $m \leq f \leq M$, and ϕ is continuous on $[m, M]$. Then $\phi \circ f \in \mathcal{R}(\alpha)$.*

Proof. Let $h = \phi \circ f$. Let $\varepsilon > 0$ be given. Since ϕ is uniformly continuous on $[m, M]$, there exists $\delta > 0$ such that $\delta < \varepsilon$, and for all $s, t \in [m, M]$,

$$|s - t| \leq \delta \implies |\phi(s) - \phi(t)| < \varepsilon.$$

Since $f \in \mathcal{R}(\alpha)$, by 15.5, there exists a partition $P = \{x_0, \dots, x_n\}$ of $[a, b]$ such that

$$U(f, \alpha; P) - L(f, \alpha; P) < \delta^2. \tag{1}$$

Let

$$\begin{aligned} M_i &= \sup_{x \in [x_{i-1}, x_i]} f(x), & M_i^* &= \sup_{x \in [x_{i-1}, x_i]} h(x), \\ m_i &= \inf_{x \in [x_{i-1}, x_i]} f(x), & m_i^* &= \inf_{x \in [x_{i-1}, x_i]} h(x). \end{aligned}$$

Divide the numbers $1, \dots, n$ into two classes:

$$A = \{i \mid M_i - m_i < \delta\},$$

$$B = \{i \mid M_i - m_i \geq \delta\}.$$

- For $i \in A$, our choice of δ shows that $M_i^* - m_i^* \leq \varepsilon$.
- For $i \in B$, $M_i^* - m_i^* \leq 2K$, where $K = \sup_{m \leq t \leq M} |\phi(t)|$.

By (1), we have

$$\delta \sum_{i \in B} \Delta\alpha_i \leq \sum_{i \in B} (M_i - m_i) \Delta\alpha_i < \delta^2$$

so that $\sum_{i \in B} \Delta\alpha_i < \delta$. It follows that

$$\begin{aligned} U(h, \alpha; P) - L(h, \alpha; P) &= \sum_{i \in A} (M_i^* - m_i^*) \Delta\alpha_i + \sum_{i \in B} (M_i^* - m_i^*) \Delta\alpha_i \\ &\leq \varepsilon (\alpha(b) - \alpha(a)) + 2K\delta \\ &< \varepsilon (\alpha(b) - \alpha(a) + 2K). \end{aligned}$$

Since ε was arbitrary, by the integrability criterion, $h \in \mathcal{R}(\alpha)$. □

§15.2 Properties of the Integral

Lemma 15.11.

(i) If $f_1, f_2 \in \mathcal{R}(\alpha)$, then $f_1 + f_2 \in \mathcal{R}(\alpha)$, and

$$\int_a^b (f_1 + f_2) d\alpha = \int_a^b f_1 d\alpha + \int_a^b f_2 d\alpha.$$

(ii) If $f \in \mathcal{R}(\alpha)$, then $cf \in \mathcal{R}(\alpha)$ for every $c \in \mathbb{R}$, and

$$\int_a^b (cf) d\alpha = c \int_a^b f d\alpha.$$

(iii) If $f_1, f_2 \in \mathcal{R}(\alpha)$ and $f_1 \leq f_2$, then

$$\int_a^b f_1 d\alpha \leq \int_a^b f_2 d\alpha.$$

(iv) If $f \in \mathcal{R}(\alpha)$ and $c \in [a, b]$, then $f \in \mathcal{R}_\alpha[a, c]$ and $f \in \mathcal{R}_\alpha[c, b]$, and

$$\int_a^b f d\alpha = \int_a^c f d\alpha + \int_c^b f d\alpha.$$

(v) If $f \in \mathcal{R}(\alpha)$ and $|f| \leq M$, then

$$\left| \int_a^b f d\alpha \right| \leq M (\alpha(b) - \alpha(a)).$$

(vi) If $f \in \mathcal{R}_{\alpha_1}[a, b]$ and $f \in \mathcal{R}_{\alpha_2}[a, b]$, then $f \in \mathcal{R}_{\alpha_1 + \alpha_2}[a, b]$, and

$$\int_a^b f d(\alpha_1 + \alpha_2) = \int_a^b f d\alpha_1 + \int_a^b f d\alpha_2;$$

if $f \in \mathcal{R}(\alpha)$ and c is a positive constant, then $f \in \mathcal{R}_{c\alpha}[a, b]$, and

$$\int_a^b f d(c\alpha) = c \int_a^b f d\alpha.$$

(vii) If $f \in \mathcal{R}(\alpha)$ and $g \in \mathcal{R}(\alpha)$, then $fg \in \mathcal{R}(\alpha)$.

Proof.

(i) If $f = f_1 + f_2$ and P is any partition of $[a, b]$, we have

$$L(f_1, \alpha; P) + L(f_2, \alpha; P) \leq L(f, \alpha; P) \leq U(f, \alpha; P) \leq U(f_1, \alpha; P) + U(f_2, \alpha; P). \quad (1)$$

If $f_1 \in \mathcal{R}(\alpha)$ and $f_2 \in \mathcal{R}(\alpha)$, let $\varepsilon > 0$ be given. There are partitions P_1 and P_2 such that

$$\begin{aligned} U(f_1, \alpha; P_1) - L(f_1, \alpha; P_1) &< \frac{\varepsilon}{2} \\ U(f_2, \alpha; P_2) - L(f_2, \alpha; P_2) &< \frac{\varepsilon}{2} \end{aligned}$$

Let P be the common refinement of P_1 and P_2 . Then (1) implies

$$U(f, \alpha; P) - L(f, \alpha; P) < \varepsilon$$

which proves that $f \in \mathcal{R}(\alpha)$.

With this same P we have

$$\begin{aligned} U(f_1, \alpha; P) &< \int_a^b f_1 \, d\alpha + \frac{\varepsilon}{2} \\ U(f_2, \alpha; P) &< \int_a^b f_2 \, d\alpha + \frac{\varepsilon}{2} \end{aligned}$$

Hence (1) implies

$$\int_a^b f \, d\alpha \leq U(f, \alpha; P) < \int_a^b f_1 \, d\alpha + \int_a^b f_2 \, d\alpha + \varepsilon.$$

Since ε was arbitrary, we conclude that

$$\int_a^b f \, d\alpha \leq \int_a^b f_1 \, d\alpha + \int_a^b f_2 \, d\alpha.$$

If we replace f_1 and f_2 in the above equation by $-f_1$ and $-f_2$, the inequality is reversed, and the equality is proved.

- (ii) The case where $c = 0$ is trivial. Given $\varepsilon > 0$, there exists P such that $U(f, \alpha; P) - L(f, \alpha; P) < \varepsilon$. If $c > 0$ write

$$U(cf, \alpha; P) = \sum_{i=1}^n cM_i\alpha_i = c \sum_{i=1}^n M_i\alpha_i = cU(f, \alpha; P).$$

Similarly,

$$L(cf, \alpha; P) = cL(f, \alpha; P).$$

Then

$$U(cf, \alpha; P) - L(cf, \alpha; P) = c(U(f, \alpha; P) - L(f, \alpha; P)) < c\varepsilon$$

and since ε is arbitrary, we are done. The case where $c < 0$ is similar. Therefore $cf \in \mathcal{R}(\alpha)$.

With this same P we have

$$U(f, \alpha; P) - \int_a^b f \, d\alpha < \varepsilon.$$

Then if $c > 0$,

$$\int_a^b cf \, d\alpha \leq U(cf, \alpha; P) = cU(f, \alpha; P) < c \int_a^b f \, d\alpha + c\varepsilon$$

so

$$\int_a^b cf \, d\alpha \leq c \int_a^b f \, d\alpha.$$

If we replace f in the above equation by $-f$, the inequality is reversed, and the equality is proved.

(iii) For every partition P , we have

$$U(f_1, \alpha; P) = \sum_{i=1}^n M_i(f_1) \Delta \alpha_i \leq \sum_{i=1}^n M_i(f_2) \Delta \alpha_i = U(f_2, \alpha; P)$$

since α is monotonically increasing on $[a, b]$.

(iv)

(v)

(vi)

(vii) Take $\phi(t) = t^2$. By 15.10, $f^2 \in R_\alpha[a, b]$ if $f \in R_\alpha[a, b]$. Write

$$fg = \frac{1}{4} \left((f+g)^2 - (f-g)^2 \right).$$

Then the desired result follows.

□

Lemma 15.12 (Triangle inequality). *Suppose $f \in \mathcal{R}(\alpha)$. Then $|f| \in \mathcal{R}(\alpha)$, and*

$$\left| \int_a^b f \, d\alpha \right| \leq \int_a^b |f| \, d\alpha.$$

Proof. Take $\phi(t) = |t|$, which is a continuous function. By 15.10, we have that $|f| = \phi \circ f \in \mathcal{R}(\alpha)$. Choose $c = \pm 1$, so that

$$c \int_a^b f \, d\alpha \geq 0.$$

Then

$$\left| \int_a^b f \, d\alpha \right| = c \int_a^b f \, d\alpha = \int_a^b cf \, d\alpha \leq \int_a^b |f| \, d\alpha,$$

since $cf \leq |f|$.

□

Example 15.13 (Heaviside step function). The *Heaviside step function* is defined by

$$H(x) = \begin{cases} 0 & (x \leq 0) \\ 1 & (x > 0) \end{cases}$$

Proposition. *Suppose f is bounded on $[a, b]$, continuous at $s \in (a, b)$. Let $\alpha(x) = H(x - s)$, then*

$$\int_a^b f \, d\alpha = f(s).$$

Proof. Consider partitions $P = \{x_0, x_1, x_2, x_3\}$, where $x_0 = a$, and $x_1 = s < x_2 < x_3 = b$. Then

$$U(f, \alpha; P) = M_2, \quad L(f, \alpha; P) = m_2.$$

Since f is continuous at s , we see that M_2 and m_2 converge to $f(s)$ as $x_2 \rightarrow s$.

□

Proposition. Suppose $c_n \geq 0$ for $n = 1, 2, \dots$, $\sum c_n$ converges, (s_n) is a sequence of distinct points in (a, b) , and

$$\alpha(x) = \sum_{n=1}^{\infty} c_n H(x - s_n).$$

Let f be continuous on $[a, b]$. Then

$$\int_a^b f \, d\alpha = \sum_{n=1}^{\infty} c_n f(s_n).$$

Proof. Since $0 \leq c_n H(x - s_n) \leq c_n$ for $n = 1, 2, \dots$ and $\sum c_n$ converges, by the comparison test, $\alpha(x) = \sum c_n H(x - s_n)$ converges for every x . Its sum $\alpha(x)$ is evidently monotonic (since each term in the sum is non-negative), and $\alpha(a) = 0$, $\alpha(b) = \sum c_n$.

Let $\varepsilon > 0$ be given. Since $\sum c_n$ converges, choose $N \in \mathbb{N}$ so that

$$\sum_{n=N+1}^{\infty} c_n < \varepsilon.$$

Let

$$\alpha_1(x) = \sum_{n=1}^N c_n H(x - s_n), \quad \alpha_2(x) = \sum_{n=N+1}^{\infty} c_n H(x - s_n).$$

By the previous result,

$$\int_a^b f \, d\alpha_1 = \sum_{n=1}^N c_n f(s_n).$$

Since $\alpha_2(b) - \alpha_2(a) < \varepsilon$,

$$\left| \int_a^b f \, d\alpha_2 \right| \leq M\varepsilon,$$

where $M = \sup |f(x)|$. Since $\alpha = \alpha_1 + \alpha_2$,

$$\int_a^b f \, d\alpha = \int_a^b f \, d\alpha_1 + \int_a^b f \, d\alpha_2$$

so it follows that

$$\left| \int_a^b f \, d\alpha - \sum_{n=1}^N c_n f(s_n) \right| \leq M\varepsilon.$$

Since ε was arbitrary, and taking $N \rightarrow \infty$, we obtain

$$\int_a^b f \, d\alpha = \sum_{n=1}^{\infty} c_n f(s_n).$$

□

In this case, we call $\alpha(x)$ a *step function*; then the integral reduces to a finite or infinite series.

The next result states that if α has an integrable derivative, then the integral reduces to an ordinary Riemann integral.

Proposition 15.14. *Assume α increases monotonically, $\alpha' \in \mathcal{R}$. Let $f: [a, b] \rightarrow \mathbb{R}$ be bounded, then $f \in \mathcal{R}(\alpha)$ if and only if $f\alpha' \in \mathcal{R}$. In that case*

$$\int_a^b f \, d\alpha = \int_a^b f(x)\alpha'(x) \, dx. \quad (15.1)$$

Proof. Let $\varepsilon > 0$ be given and apply 15.5 to α' : There exists a partition $P = \{x_0, \dots, x_n\}$ of $[a, b]$ such that

$$U(\alpha'; P) - L(\alpha'; P) < \varepsilon. \quad (1)$$

By the mean value theorem, there exist points $t_i \in [x_{i-1}, x_i]$ such that

$$\Delta\alpha_i = \alpha'(t_i)\Delta x_i \quad (i = 1, \dots, n).$$

If $s_i \in [x_{i-1}, x_i]$, then by 15.6,

$$\sum_{i=1}^n |\alpha'(s_i) - \alpha'(t_i)| \Delta x_i < \varepsilon. \quad (2)$$

Let $M = \sup |f(x)|$. Since

$$\sum_{i=1}^n f(s_i)\Delta\alpha_i = \sum_{i=1}^n f(s_i)\alpha'(t_i)\Delta x_i$$

it follows from (2) that

$$\begin{aligned} \left| \sum_{i=1}^n f(s_i)\Delta\alpha_i - \sum_{i=1}^n f(s_i)\alpha'(s_i)\Delta x_i \right| &= \left| \sum_{i=1}^n f(s_i) (\alpha'(t_i) - \alpha'(s_i)) \Delta x_i \right| \\ &\leq \sum_{i=1}^n |f(s_i)| |\alpha'(t_i) - \alpha'(s_i)| \Delta x_i \\ &= \sum_{i=1}^n |f(s_i)| |\alpha'(t_i) - \alpha'(s_i)| \Delta x_i \\ &\leq M \sum_{i=1}^n |\alpha'(t_i) - \alpha'(s_i)| \Delta x_i \\ &\leq M\varepsilon. \end{aligned} \quad (3)$$

In particular, for all choices of $s_i \in [x_{i-1}, x_i]$,

$$\sum_{i=1}^n f(s_i)\Delta\alpha_i \leq U(f\alpha'; P) + M\varepsilon$$

so taking sup for $f(s_i)$ gives

$$U(f, \alpha; P) \leq U(f\alpha'; P) + M\varepsilon.$$

The same argument leads from (3) to

$$U(f\alpha'; P) \leq U(f, \alpha; P) + M\varepsilon.$$

Hence

$$|U(f, \alpha; P) - U(f\alpha'; P)| \leq M\varepsilon. \quad (4)$$

Since (1) holds true for any refinement of P , hence (4) also remains true. We conclude that

$$\left| \int_a^{\bar{b}} f \, d\alpha - \int_a^{\bar{b}} f(x)\alpha'(x) \, dx \right| \leq M\varepsilon.$$

But ε is arbitrary. Hence

$$\int_a^{\bar{b}} f \, d\alpha = \int_a^{\bar{b}} f(x)\alpha'(x) \, dx$$

for any bounded f . The equality of the lower integrals follows from

$$\begin{aligned} \int_a^{\bar{b}} -f \, d\alpha &= \int_a^{\bar{b}} -f\alpha' \, dx \\ - \int_a^{\bar{b}} f \, d\alpha &= - \int_a^{\bar{b}} f\alpha' \, dx \\ \int_a^{\bar{b}} f \, d\alpha &= \int_a^{\bar{b}} f(x)\alpha'(x) \, dx \end{aligned}$$

Therefore the theorem follows. □

Proposition 15.15 (Change of variables). *Suppose $\phi: [A, B] \rightarrow [a, b]$ is strictly increasing and continuous. Suppose α is monotonically increasing on $[a, b]$, $f \in \mathcal{R}(\alpha)$. Define β and g on $[A, B]$ by*

$$\beta(y) = \alpha(\phi(y)), \quad g(y) = f(\phi(y)).$$

Then $g \in \mathcal{R}(\beta)$, and

$$\int_A^B g \, d\beta = \int_a^b f \, d\alpha. \tag{15.2}$$

Proof. To each partition $P = \{x_0, \dots, x_n\}$ of $[a, b]$ corresponds a partition $Q = \{y_0, \dots, y_n\}$ of $[A, B]$, where

$$x_i = \phi(y_i) \quad (i = 1, \dots, n).$$

All partitions of $[A, B]$ are obtained in this way. Since the values taken by f on $[x_{i-1}, x_i]$ are exactly the same as those taken by g on $[y_{i-1}, y_i]$, we see that

$$\begin{aligned} U(g, \beta; Q) &= U(f, \alpha; P), \\ L(g, \beta; Q) &= L(f, \alpha; P). \end{aligned} \tag{1}$$

Since $f \in \mathcal{R}(\alpha)$, P can be chosen so that both $U(f, \alpha; P)$ and $L(f, \alpha; P)$ are close to $\int f \, d\alpha$. Hence (1), combined with 15.5, shows that $g \in \mathcal{R}_\beta[A, B]$ and

$$\int_A^B g \, d\beta = \int_a^b f \, d\alpha.$$

□

Note the following special case: Take $\alpha(x) = x$. Then $\beta = \phi$. Assume $\phi' \in \mathcal{R}$. Applying 15.14 to the LHS of

$$\int_A^B g \, d\beta = \int_a^b f \, d\alpha,$$

we obtain

$$\int_a^b f(x) \, dx = \int_A^B f(\phi(y)) \phi'(y) \, dy.$$

§15.3 Integration and Differentiation

We shall show that integration and differentiation are, in a certain sense, inverse operations.

Theorem 15.16. *Suppose $f \in \mathcal{R}(\alpha)$. For $a \leq x \leq b$, let the cumulative function be*

$$F(x) = \int_a^x f(t) dt.$$

Then F is continuous on $[a, b]$; furthermore, if f is continuous at $x_0 \in [a, b]$, then F is differentiable at x_0 , and

$$F'(x_0) = f(x_0).$$

Proof. Suppose $f \in \mathcal{R}(\alpha)$. Since f is bounded, let $|f(t)| \leq M$ for $t \in [a, b]$. If $a \leq x < y \leq b$, then

$$\begin{aligned} |F(y) - F(x)| &= \left| \int_a^y f(t) dt - \int_a^x f(t) dt \right| \\ &= \left| \int_x^y f(t) dt \right| \\ &\leq \int_x^y |f(t)| dt \\ &\leq M(y - x). \end{aligned}$$

Hence F is Lipschitz continuous, so F is uniformly continuous on $[a, b]$.

Now suppose f is continuous at x_0 . Fix $\varepsilon > 0$, choose $\delta > 0$ such that for $a \leq t \leq b$,

$$|t - x_0| < \delta \implies |f(t) - f(x_0)| < \varepsilon.$$

Hence, if s, t are such that

$$x_0 - \delta < s \leq x_0 \leq t < x_0 + \delta \quad \text{and} \quad a \leq x < t \leq b,$$

we have, by 15.11(v),

$$\begin{aligned} \left| \frac{F(t) - F(s)}{t - s} - f(x_0) \right| &= \left| \frac{\int_a^s f(u) du - \int_a^s f(u) du}{t - s} - f(x_0) \right| \\ &= \left| \frac{1}{t - s} \int_s^t (f(u) - f(x_0)) du \right| \\ &= \frac{1}{t - s} \left| \int_s^t (f(u) - f(x_0)) du \right| \\ &\leq \frac{1}{t - s} \int_s^t |f(u) - f(x_0)| du \\ &< \frac{1}{t - s} \varepsilon(t - s) = \varepsilon \end{aligned}$$

so it follows that $F'(x_0) = f(x_0)$. □

Theorem 15.17 (Fundamental theorem of calculus). *Suppose $f \in \mathcal{R}(\alpha)$, and there exists a differentiable function F on $[a, b]$ such that $F' = f$. Then*

$$\int_a^b f(x) \, dx = F(b) - F(a). \quad (15.3)$$

Proof. Let $\varepsilon > 0$ be given. Choose a partition $P = \{x_0, \dots, x_n\}$ of $[a, b]$ such that $U(f; P) - L(f; P) < \varepsilon$. By the mean value theorem, there exist $t_i \in [x_{i-1}, x_i]$ such that

$$\begin{aligned} F(x_i) - F(x_{i-1}) &= F'(t_i)\Delta x_i \\ &= f(t_i)\Delta x_i. \end{aligned}$$

Thus

$$\sum_{i=1}^n f(t_i)\Delta x_i = F(b) - F(a).$$

Then by 15.6,

$$\left| F(b) - F(a) - \int_a^b f(x) \, dx \right| = \left| \sum_{i=1}^n f(t_i)\Delta x_i - \int_a^b f(x) \, dx \right| < \varepsilon.$$

Since this holds for all $\varepsilon > 0$, the proof is complete. \square

Lemma 15.18 (Integration by parts). *Suppose F and G are differentiable on $[a, b]$, $F' = f \in \mathcal{R}$ and $G' = g \in \mathcal{R}$. Then*

$$\int_a^b F(x)g(x) \, dx = F(b)G(b) - F(a)G(a) - \int_a^b f(x)G(x) \, dx. \quad (15.4)$$

Proof. Let $H(x) = F(x)G(x)$. Then apply the fundamental theorem of calculus to H and its derivative. \square

§15.4 Integration of Vector-valued Functions

Let $f_1, \dots, f_k: [a, b] \rightarrow \mathbb{R}$, and let $\mathbf{f} = (f_1, \dots, f_k)$ where $\mathbf{f}: [a, b] \rightarrow \mathbb{R}^k$. We say that $\mathbf{f} \in \mathcal{R}(\alpha)$ if $f_1, \dots, f_k \in \mathcal{R}(\alpha)$. If this is the case, we define

$$\int_a^b \mathbf{f} \, d\alpha := \left(\int_a^b f_1 \, d\alpha, \dots, \int_a^b f_k \, d\alpha \right).$$

In other words, we “integrate componentwise”, so that $\int \mathbf{f} \, d\alpha$ is the point in \mathbb{R}^k whose i -th coordinate is $\int f_i \, d\alpha$.

It is clear that parts (a), (c), and (e) of Theorem 6.12 are valid for these vector-valued integrals; we simply apply the earlier results to each coordinate. The same is true of Theorems 6.17, 6.20, and 6.21. To illustrate, we state the analogue of the fundamental theorem of calculus.

Theorem 15.19. *If $\mathbf{f}, \mathbf{F}: [a, b] \rightarrow \mathbb{R}^k$, $\mathbf{f} \in \mathcal{R}(\alpha)$, and $\mathbf{F}' = \mathbf{f}$. Then*

$$\int_a^b \mathbf{f}(t) \, dt = \mathbf{F}(b) - \mathbf{F}(a). \quad (15.5)$$

The analogue of Theorem 6.13(b) offers some new features, however, at least in its proof.

Lemma 15.20 (Triangle inequality). *Let $\mathbf{f}: [a, b] \rightarrow \mathbb{R}^k$, $\mathbf{f} \in \mathcal{R}(\alpha)$ where α is monotonically increasing on $[a, b]$. Then $\|\mathbf{f}\| \in \mathcal{R}(\alpha)$, and*

$$\left\| \int_a^b \mathbf{f} \, d\alpha \right\| \leq \int_a^b \|\mathbf{f}\| \, d\alpha.$$

Proof. If f_1, \dots, f_k are the components of \mathbf{f} , then

$$\|\mathbf{f}\| = \left(f_1^2 + \dots + f_k^2 \right)^{1/2}.$$

By 15.10, each of the functions $f_i^2 \in \mathcal{R}(\alpha)$, so their sum $f_1^2 + \dots + f_k^2 \in \mathcal{R}(\alpha)$.

Since x^2 is a continuous function of x , Theorem 4.17 shows that the square-root function is continuous on $[0, M]$, for every real M . If we apply Theorem 6.11 once more, (41) shows that $\|\mathbf{f}\| \in \mathcal{R}(\alpha)$.

Let $\mathbf{y} = (y_1, \dots, y_k)$, where $y_i = \int f_i \, d\alpha$. Then we have $\mathbf{y} = \int \mathbf{f} \, d\alpha$, and

$$\|\mathbf{y}\|^2 = \sum_{i=1}^k y_i^2 = \sum_{i=1}^k \left(y_i \int f_i \, d\alpha \right) = \int \left(\sum_{i=1}^k y_i f_i \right) d\alpha.$$

By the Cauchy–Schwarz inequality,

$$\sum_{i=1}^k y_i f_i(t) \leq \|\mathbf{y}\| \|\mathbf{f}(t)\| \quad (a \leq t \leq b);$$

hence Theorem 6.12(b) implies

$$\|\mathbf{f}\|^2 \leq \|\mathbf{y}\| \int \|\mathbf{f}\| \, d\alpha.$$

If $\mathbf{y} = \mathbf{0}$, (40) is trivial. If $\mathbf{y} \neq \mathbf{0}$, division of (43) by $\|\mathbf{y}\|$ gives (40).

□

to do

§15.5 Rectifiable Curves

Definition 15.21 (Curve). A **curve** in \mathbb{R}^k is a continuous mapping $\gamma: [a, b] \rightarrow \mathbb{R}^k$. If γ is bijective, γ is called an **arc**. If $\gamma(a) = \gamma(b)$, γ is said to be a **closed curve**.

The case $k = 2$ (i.e., the case of plane curves) is of considerable importance in the study of analytic functions of a complex variable.

Remark. Note that we define a curve to be a mapping, not a point set. Of course, with each curve γ in \mathbb{R}^k there is associated a subset of \mathbb{R}^k , namely the range of γ , but different curves may have the same range.

For each partition $P = \{x_0, \dots, x_n\}$ of $[a, b]$ and each curve γ on $[a, b]$, define

$$\Lambda(\gamma; P) := \sum_{i=1}^n |\gamma(x_i) - \gamma(x_{i-1})|.$$

The i -th term in this sum is the distance (in \mathbb{R}^k) between the points $\gamma(x_{i-1})$ and $\gamma(x_i)$. Hence $\Lambda(\gamma; P)$ is the length of a polygonal path with vertices at $\gamma(x_0), \gamma(x_1), \dots, \gamma(x_n)$, in this order. As our partition becomes finer and finer, this polygon approaches the range of γ more and more closely.

Definition 15.22. The **total variation** (or **length**) of γ is

$$\Lambda(\gamma) := \sup_{P \in \mathcal{P}[a, b]} \Lambda(\gamma; P).$$

We say γ is **rectifiable** if $\Lambda(\gamma) < \infty$.

The next result gives a formula for calculating the length of a rectifiable curve that is continuously differentiable.

Proposition 15.23. *If γ is a continuously differentiable curve on $[a, b]$, then γ is rectifiable, and*

$$\Lambda(\gamma) = \int_a^b |\gamma'(t)| dt. \quad (15.6)$$

Proof. If $a \leq x_{i-1} < x_i \leq b$, then

$$|\gamma(x_i) - \gamma(x_{i-1})| = \left| \int_{x_{i-1}}^{x_i} \gamma'(t) dt \right| \leq \int_{x_{i-1}}^{x_i} |\gamma'(t)| dt.$$

Hence, for every partition P of $[a, b]$, taking the sum on both sides gives

$$\Lambda(\gamma; P) \leq \int_a^b |\gamma'(t)| dt$$

and taking sup gives

$$\Lambda(\gamma) \leq \int_a^b |\gamma'(t)| dt.$$

We now prove the opposite inequality. Since γ' is (continuous and thus) uniformly continuous on $[a, b]$, fix

insert
figure

$\varepsilon > 0$, there exists $\delta > 0$ such that

$$|s - t| < \delta \implies |\gamma'(s) - \gamma'(t)| < \varepsilon.$$

Let $P = \{x_0, \dots, x_n\}$ be a partition of $[a, b]$, with $\Delta x_i < \delta$ for all i . If $t \in [x_{i-1}, x_i]$, it follows that

$$|\gamma'(t)| \leq |\gamma'(x_i)| + \varepsilon.$$

Hence

$$\begin{aligned} \int_{x_{i-1}}^{x_i} |\gamma'(t)| dt &\leq |\gamma'(x_i)|\Delta x_i + \varepsilon\Delta x_i \\ &= \left| \int_{x_{i-1}}^{x_i} (\gamma'(t) + \gamma'(x_i) - \gamma'(t)) dt \right| + \varepsilon\Delta x_i \\ &\leq \left| \int_{x_{i-1}}^{x_i} \gamma'(t) dt \right| + \left| \int_{x_{i-1}}^{x_i} (\gamma'(x_i) - \gamma'(t)) dt \right| + \varepsilon\Delta x_i \\ &\leq |\gamma(x_i) - \gamma(x_{i-1})| + 2\varepsilon\Delta x_i. \end{aligned}$$

If we add these inequalities, we obtain

$$\begin{aligned} \int_a^b |\gamma'(t)| dt &\leq \Lambda(\gamma; P) + 2\varepsilon(b - a) \\ &\leq \Lambda(\gamma) + 2\varepsilon(b - a). \end{aligned}$$

Since ε was arbitrary, we must have

$$\int_a^b |\gamma'(t)| \leq \Lambda(\gamma).$$

This completes the proof. □

Exercises

16 Sequences and Series of Functions

Suppose $f_n : E \subset X \rightarrow Y$ is a sequence of functions. In some cases, we shall restrict ourselves to complex-valued functions (take $Y = \mathbb{C}$).

§16.1 Pointwise Convergence

A natural extension of convergence of sequences of numbers to sequences of functions is to fix a point $x \in E$, and consider the behaviour of the sequence $(f_n(x))$.

Definition 16.1 (Pointwise convergence). Suppose (f_n) is a sequence of functions, and $(f_n(x))$ converges for every $x \in E$. We say (f_n) **converges pointwise** to f on E , denoted by $f_n \rightarrow f$, if

$$f(x) = \lim_{n \rightarrow \infty} f_n(x) \quad (\forall x \in E).$$

That is, for all $x \in E$,

$$\forall \varepsilon > 0, \quad \exists N \in \mathbb{N}, \quad \forall n \geq N, \quad d(f_n(x) - f(x)) < \varepsilon.$$

f is called the *limit* (or *limit function*) of (f_n) .

Similarly, if $\sum f_n(x)$ converges for every $x \in E$, and if we define

$$f(x) = \sum_{n=1}^{\infty} f_n(x) \quad (\forall x \in E)$$

the function f is called the *sum of the series* $\sum f_n$.

Example 16.2. The sequence of functions $f_n(x) = \frac{x}{n}$ converges pointwise to the zero function $f(x) = 0$.

The main problem which arises is to determine whether important properties of functions are preserved by pointwise convergence. For instance, if f_n are continuous, or differentiable, or integrable, is the same true of the limit function? What are the relations between f'_n and f' , say, or between $\int f_n$ and $\int f$?

Example 16.3 (Continuity). For $0 < x < 1$, the sequence of functions $f_n(x) = x^n$ converges pointwise to the function

$$f(x) = \begin{cases} 1 & (x = 1) \\ 0 & (0 \leq x < 1) \end{cases}$$

Evidently f_n are continuous, but f is discontinuous. Hence

$$\lim_{x \rightarrow x_0} \lim_{n \rightarrow \infty} f_n(x) \neq \lim_{n \rightarrow \infty} \lim_{x \rightarrow x_0} f_n(x).$$

Example 16.4 (Differentiability). For $x \in \mathbb{R}$, let

$$f_n(x) = \frac{\sin nx}{\sqrt{n}} \quad (n = 1, 2, \dots)$$

so

$$f(x) = \lim_{n \rightarrow \infty} f_n(x) = 0.$$

Then $f'(x) = 0$, and

$$f'_n(x) = \sqrt{n} \cos nx,$$

so (f'_n) does not converge to f' .

This shows that the limit of the derivative does not equal the derivative of the limit.

Example 16.5 (Integrability). Let

$$f_n(x) = \chi_{[n, n+1]}(x),$$

Then $\int_{\mathbb{R}} f_n(x) \, dx = 1$, so

$$\lim_{n \rightarrow \infty} \int_{\mathbb{R}} f_n(x) \, dx = 1.$$

However

$$\int_{\mathbb{R}} \lim_{n \rightarrow \infty} f_n(x) \, dx = \int_{\mathbb{R}} 0 \, dx = 0.$$

This shows that the limit of the integral does not equal the integral of the limit. Thus we may not switch the order of limits.

Pointwise convergence does not preserve many nice properties of functions. Hence, we need a stronger notion of convergence for sequences and series of functions.

§16.2 Uniform Convergence

Definition 16.6 (Uniform convergence). We say (f_n) *converges uniformly* to f over E , denoted by $f_n \rightrightarrows f$, if

$$\forall \varepsilon > 0, \quad \exists N \in \mathbb{N}, \quad \forall x \in E, \quad \forall n \geq N, \quad d(f_n(x) - f(x)) < \varepsilon.$$

Similarly, a series of functions $\sum f_n(x)$ converges uniformly on E if the sequence of partial sums (s_n) defined by

$$s_n(x) = \sum_{k=1}^n f_k(x)$$

converges uniformly on E .

Intuitively, uniform convergence can be visualised as the sequence of functions (f_n) eventually contained in an ε -tube around f , for sufficiently large n .

Remark. Uniform convergence is stronger than pointwise convergence, since N is uniform (or “fixed”) for all $x \in E$; for pointwise convergence, the choice of N is determined by x .

Uniform convergence implies pointwise convergence, but not the other way around.

Example 16.7. Consider the sequence of functions $f_n(x) = x^n$ defined on $(0, 1)$. Then $f_n \rightarrow 0$. But $f_n \not\rightrightarrows 0$.

Proof. □

We shall restrict our focus to sequences of complex-valued functions defined on $E \subset X$, for the remaining of the chapter.

The Cauchy criterion for uniform convergence is as follows.

Lemma 16.8 (Cauchy criterion). $f_n \rightrightarrows f$ on E if and only if

$$\forall \varepsilon > 0, \quad \exists N \in \mathbb{N}, \quad \forall x \in E, \quad \forall n, m \geq N, \quad |f_n(x) - f_m(x)| < \varepsilon.$$

Proof.

\Rightarrow Suppose $f_n \rightrightarrows f$ on E . Let $\varepsilon > 0$ be given. Then there exists $N \in \mathbb{N}$ such that for all $x \in E$, for all $n \geq N$,

$$|f_n(x) - f(x)| < \frac{\varepsilon}{2}.$$

Then for all $n, m \geq N$,

$$\begin{aligned} |f_n(x) - f_m(x)| &= |(f_n(x) - f(x)) + (f(x) - f_m(x))| \\ &\leq |f_n(x) - f(x)| + |f_m(x) - f(x)| \\ &< \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon. \end{aligned}$$

\Leftarrow Suppose the Cauchy criterion holds for (f_n) ; that is,

$$\forall \varepsilon > 0, \quad \exists N \in \mathbb{N}, \quad \forall x \in E, \quad \forall n, m \geq N, \quad |f_n(x) - f_m(x)| < \varepsilon.$$

insert
figure

Then for every $x \in E$, the sequence $(f_n(x))$ is a Cauchy sequence and thus converges to a limit $f(x)$. Hence by definition, $f_n \rightarrow f$ on E . We are left to prove that the convergence is uniform.

Let $\varepsilon > 0$ be given. There exists $N \in \mathbb{N}$ such that for all $n, m \geq N$ and for all $x \in E$,

$$|f_n(x) - f_m(x)| < \varepsilon.$$

Fix n , and let $m \rightarrow \infty$. Since $\lim_{m \rightarrow \infty} f_m(x) = f(x)$, thus for all $n \geq N$ and for all $x \in E$,

$$|f_n(x) - f(x)| < \varepsilon,$$

which completes the proof. □

Definition 16.9. If $f \in \mathcal{C}(X, \mathbb{C})$, we define the *supremum norm* of f as

$$\|f\| := \sup_{x \in X} |f(x)|.$$

Lemma 16.10. $\|f\|$ gives a norm on $\mathcal{C}(X, \mathbb{C})$. Then $\mathcal{C}(X, \mathbb{C})$ is a metric space, with metric $d(f, g) = \|f - g\|$.

Proof. Check that $\|f\|$ satisfies the conditions for a norm:

(i) $|f(x)| \geq 0$ for all $x \in X$, so $\|f\| \geq 0$. It is clear that $\|f\| = 0$ if and only if $f(x) = 0$ for every $x \in X$, that is, only if $f = 0$.

(ii) For all $\lambda \in \mathbb{C}$,

$$\|\lambda f\| = \sup_{x \in X} |\lambda f(x)| = |\lambda| \sup_{x \in X} |f(x)| = |\lambda| \|f\|.$$

(iii) If $h = f + g$, then for all $x \in X$,

$$|h(x)| \leq |f(x)| + |g(x)| \leq \|f\| + \|g\|.$$

Hence taking sup on the left gives $\|f + g\| \leq \|f\| + \|g\|$.

Check conditions for metric space. □

The following result provides another way to determine uniform convergence.

Lemma 16.11. $f_n \rightrightarrows f$ on E if and only if $f_n \rightarrow f$ on E with respect to the metric of $\mathcal{C}(E, \mathbb{C})$.

Proof.

$$\begin{aligned}
 f_n \rightarrow f &\iff \lim_{n \rightarrow \infty} \|f_n - f\| = 0 \\
 &\iff \lim_{n \rightarrow \infty} \left(\sup_{x \in E} |f_n(x) - f(x)| \right) = 0 \\
 &\iff \forall \varepsilon > 0, \exists N \in \mathbb{N}, \forall n \geq N, \sup_{x \in E} |f_n(x) - f(x)| < \varepsilon \\
 &\iff \forall \varepsilon > 0, \exists N \in \mathbb{N}, \forall n \geq N, \forall x \in E, |f_n(x) - f(x)| < \varepsilon
 \end{aligned}$$

which precisely means that $f_n \rightrightarrows f$ on E , by definition.

Note that for the last step, the $\boxed{\Leftarrow}$ direction is tricky, since the limit can equal ε , so we take $\frac{\varepsilon}{2}$ instead. \square

For series, there is a very convenient test for uniform convergence, due to Weierstrass.

Lemma 16.12 (Weierstrass M-test). *Suppose (f_n) is a sequence of complex-valued functions defined on E , and*

$$|f_n(x)| \leq M_n \quad (n = 1, 2, \dots, x \in E)$$

If $\sum M_n$ converges, then $\sum f_n$ converges uniformly on E .

Proof. Suppose $\sum M_n$ converges. Let $\varepsilon > 0$ be given, the partial sums of $\sum M_n$ form a Cauchy sequence, so there exists $N \in \mathbb{N}$ such that for all $n \geq m \geq N$,

$$\sum_{k=m}^n M_k < \varepsilon.$$

Then considering the partial sums of the series of functions,

$$\left| \sum_{k=m}^n f_k(x) \right| \leq \sum_{k=m}^n |f_k(x)| \leq \sum_{k=m}^n M_k < \varepsilon.$$

By the Cauchy criterion (16.8), we are done. \square

Example 16.13.

- The series $\sum_{n=1}^{\infty} \frac{\sin nx}{n^2}$ converges uniformly on \mathbb{R} . (Note: this is a Fourier series, we'll see more of these later). That is because

$$\left| \frac{\sin nx}{n^2} \right| \leq \frac{1}{n^2} \quad \text{and} \quad \sum_{n=1}^{\infty} \frac{1}{n^2} \text{ converges.}$$

- The series $\sum_{n=0}^{\infty} \frac{x^n}{n!}$ converges uniformly on any bounded interval. For example take the interval $[-r, r] \subset \mathbb{R}$,

$$\left| \frac{x^n}{n!} \right| \leq \frac{r^n}{n!} \quad \text{and} \quad \sum_{n=1}^{\infty} \frac{r^n}{n!} \text{ converges by the ratio test.}$$

§16.3 Properties of Uniform Convergence

We now consider properties preserved by uniform convergence.

Uniform Convergence and Continuity

We prove a more general result.

Proposition 16.14. *Suppose $f_n \rightrightarrows f$ on E . Let $x \in X$ be a limit point of E , and suppose that*

$$\lim_{t \rightarrow x} f_n(t) = A_n \quad (n = 1, 2, \dots).$$

Then (A_n) converges, and $\lim_{t \rightarrow x} f(t) = \lim_{n \rightarrow \infty} A_n$.

In other words, the conclusion is that

$$\lim_{t \rightarrow x} \lim_{n \rightarrow \infty} f_n(t) = \lim_{n \rightarrow \infty} \lim_{t \rightarrow x} f_n(t).$$

Proof.

1. We first show that (A_n) converges. Since (f_n) uniformly converges on E , by the Cauchy criterion (16.8), fix $\varepsilon > 0$, there exists $N \in \mathbb{N}$ such that for all $n, m \geq N, t \in E$,

$$|f_n(t) - f_m(t)| < \varepsilon.$$

Letting $t \rightarrow x$, since $\lim_{t \rightarrow x} f_n(t) = A_n$, we have that for all $n, m \geq N$,

$$|A_n - A_m| < \varepsilon.$$

Thus (A_n) is a Cauchy sequence and therefore converges, say to A .

2. Next we will show that $\lim_{t \rightarrow x} f(t) = A$.

Idea. We want to bound the term $|f(t) - A|$, using terms of known values.

Write

$$|f(t) - A| \leq |f(t) - f_n(t)| + |f_n(t) - A_n| + |A_n - A|. \quad (1)$$

By the uniform convergence of (f_n) , there exists $N_1 \in \mathbb{N}$ such that for all $n \geq N_1$,

$$|f(t) - f_n(t)| < \frac{\varepsilon}{3} \quad (t \in E).$$

By the convergence of (A_n) , there exists $N_2 \in \mathbb{N}$ such that for all $n \geq N_2$,

$$|A_n - A| < \frac{\varepsilon}{3}.$$

Choose $N = \max\{N_1, N_2\}$ such that the above two inequalities hold simultaneously. Then for this n ,

since $\lim_{t \rightarrow x} f_n(t) = A_n$, we choose an open ball B of x such that if $t \in B \cap E$, $t \neq x$, then

$$|f_n(t) - A_n| < \frac{\varepsilon}{3}.$$

Substituting the above inequalities into (1) gives

$$|f(t) - A| < \frac{\varepsilon}{3} + \frac{\varepsilon}{3} + \frac{\varepsilon}{3} = \varepsilon.$$

provided $t \in B \cap E$, $t \neq x$. This is equivalent to $\lim_{t \rightarrow x} f(t) = A$.

□

An immediate important corollary is that uniform convergence preserves continuity.

Corollary 16.15. *Suppose (f_n) are continuous on E , and $f_n \rightrightarrows f$ on E . Then f is continuous on E .*

Proof. By continuity of f_n ,

$$\lim_{t \rightarrow x} f_n(t) = f_n(x).$$

Then

$$\lim_{t \rightarrow x} f(t) = \lim_{t \rightarrow x} \left(\lim_{n \rightarrow \infty} f_n(t) \right) = \lim_{n \rightarrow \infty} \left(\lim_{t \rightarrow x} f_n(t) \right) = \lim_{n \rightarrow \infty} f_n(x) = f(x),$$

which precisely means that f is continuous on E . □

Remark. The converse is not true; for instance, the sequence of functions $f_n : (0, 1) \rightarrow \mathbb{R}$ defined by $f_n(x) = x^n$ converges to the zero function, which is continuous, but the convergence is not uniform.

Let us see that we can have extra conditions such that the converse of the previous result is true.

Proposition 16.16 (Dini's theorem). *Suppose K is compact, and (f_n) is a sequence of continuous functions on K , $f_n \rightarrow f$ on K , and (f_n) is monotonically decreasing:*

$$f_n(x) \geq f_{n+1}(x) \quad (n = 1, 2, \dots).$$

Then $f_n \rightrightarrows f$ on K .

Remark. The compactness in the hypotheses is necessary; for instance, on $(0, 1)$ define $f_n(x) = \frac{1}{nx + 1}$. Then $f_n(x) \rightarrow 0$ monotonically in $(0, 1)$, but the convergence is not uniform.

Proof. Let $g_n = f_n - f$. Then g_n is continuous, $g_n \rightarrow 0$, and $g_n \geq g_{n+1} \geq 0$. We have to prove that $g_n \rightrightarrows 0$ on K .

Let $\varepsilon > 0$ be given. For $n = 1, 2, \dots$, let

$$K_n = \{x \in K \mid g_n(x) \geq \varepsilon\}.$$

Since g_n is continuous, and $\{g_n(x) \mid g_n(x) \geq \varepsilon\}$ is closed, by 13.13, its pre-image K_n is closed. Since K_n is a closed subset of a compact set K , by 11.39, K_n is compact.

Since $g_n \geq g_{n+1}$, we have $K_n \supset K_{n+1}$. Fix $x \in K$. Since $g_n(x) \rightarrow 0$, we see that $x \notin K_n$ if n is sufficiently large. Thus $x \notin \bigcap_{n=1}^{\infty} K_n$. In other words, $\bigcap_{n=1}^{\infty} K_n = \emptyset$. Hence $K_N = \emptyset$ for some N (by the converse of

Cantor's intersection theorem). It follows that for all $x \in K$ and for all $n \geq N$,

$$0 \leq g_n(x) < \varepsilon.$$

Therefore $g_n \rightrightarrows 0$ on K , as desired. \square

Lemma 16.17. $\mathcal{C}(X, \mathbb{C})$ is a complete metric space.

Proof. Let (f_n) be a Cauchy sequence in $\mathcal{C}(X, \mathbb{C})$. Then fix $\varepsilon > 0$, there exists $N \in \mathbb{N}$ such that for all $n, m \geq N$,

$$\|f_n - f_m\| < \varepsilon.$$

By the Cauchy criterion (16.8), $f_n \rightrightarrows f$ for some $f : X \rightarrow \mathbb{C}$. We now need to show that $f \in \mathcal{C}(X, \mathbb{C})$; that is, f is continuous and bounded.

- f is continuous by 16.15.
- f is bounded, since there is an n such that $|f(x) - f_n(x)| < 1$ for all $x \in X$, and f_n is bounded.

Hence $f \in \mathcal{C}(X, \mathbb{C})$, and since $f_n \rightrightarrows f$ on X , we have $\|f - f_n\| \rightarrow 0$ as $n \rightarrow \infty$. \square

Uniform Convergence and Integration

The next result states that the limit and integral can be interchanged.

Proposition 16.18. *Suppose (f_n) are defined over $[a, b]$ and $f_n \in \mathcal{R}(\alpha)$. If $f_n \rightrightarrows f$ on $[a, b]$, then $f \in \mathcal{R}(\alpha)$, and*

$$\lim_{n \rightarrow \infty} \int_a^b f_n \, d\alpha = \int_a^b f \, d\alpha.$$

Proof. It suffices to prove this for real-valued f_n . Let

to do

$$\varepsilon_n = \sup_{x \in [a, b]} |f_n(x) - f(x)|.$$

Then

$$f_n - \varepsilon_n \leq f \leq f_n + \varepsilon_n,$$

so that the upper and lower integrals of f (see Definition 6.2) satisfy

$$\int_a^b (f_n - \varepsilon_n) \, d\alpha \leq \int_a^b f \, d\alpha \leq \int_a^b f \, d\alpha \leq \int_a^b (f_n + \varepsilon_n) \, d\alpha.$$

Hence

$$0 \leq \int_a^b f \, d\alpha - \int_a^b f \, d\alpha \leq 2\varepsilon_n[\alpha(b) - \alpha(a)].$$

Since $\varepsilon_n \rightarrow 0$ as $n \rightarrow \infty$ (Theorem 7.9), the upper and lower integrals of f are equal. Thus $f \in \mathcal{R}(\alpha)$.

Another application of (25) now yields

$$\left| \int_a^b f \, d\alpha - \int_a^b f_n \, d\alpha \right| \leq \varepsilon_n[\alpha(b) - \alpha(a)].$$

This implies

$$\lim_{n \rightarrow \infty} \int_a^b f_n \, d\alpha = \int_a^b f \, d\alpha.$$

□

Corollary 16.19. *Suppose $f_n \in \mathcal{R}(\alpha)$ and*

$$f(x) = \sum_{n=1}^{\infty} f_n(x)$$

converges uniformly on $[a, b]$. Then

$$\int_a^b f \, d\alpha = \sum_{n=1}^{\infty} \int_a^b f_n \, d\alpha.$$

In other words, we can swap the integral and sum, such that the series may be integrated term by term.

Proof. Consider the sequence of partial sums

$$f_n(x) = \sum_{k=1}^n f_k(x) \quad (n = 1, 2, \dots).$$

It follows $f_n \in \mathcal{R}(\alpha)$ and $f_n \rightrightarrows f$. Apply above theorem to (f_n) and the conclusion follows. \square

Example 16.20. Let us show how to integrate a Fourier series:

$$\int_0^x \sum_{n=1}^{\infty} \frac{\cos nt}{n^2} dt = \sum_{n=1}^{\infty} \int_0^x \frac{\cos nt}{n^2} dt = \sum_{n=1}^{\infty} \frac{\sin nx}{n^3}.$$

Uniform Convergence and Differentiation

The next result shows that the process of limit and differentiation can be interchanged.

Proposition 16.21. *Suppose (f_n) are differentiable on $[a, b]$, and $(f_n(x_0))$ converges for some $x_0 \in [a, b]$. If f'_n converges uniformly on $[a, b]$, then there exists a differentiable f such that $f_n \rightrightarrows f$ on $[a, b]$, and*

$$f'(x) = \lim_{n \rightarrow \infty} f'_n(x) \quad (a \leq x \leq b).$$

Proof. Let $\varepsilon > 0$ be given. Since $(f_n(x_0))$ converges, $(f_n(x_0))$ is a Cauchy sequence, so there exists $N \in \mathbb{N}$ such that for all $n, m \geq N$,

$$|f_n(x_0) - f_m(x_0)| < \frac{\varepsilon}{2}.$$

Since (f'_n) converges uniformly on $[a, b]$, by 16.8,

$$|f'_n(x) - f'_m(x)| < \frac{\varepsilon}{2(b-a)} \quad (a \leq x \leq b).$$

Now apply the mean value theorem to the function $f_n - f_m$: for $x_0, x \in [a, b]$, there exists t between x_0 and x such that

$$(f_n - f_m)(x_0) - (f_n - f_m)(x) = (f_n - f_m)'(t)(x_0 - x)$$

and thus if $n, m \geq N$, then

$$\begin{aligned} \left| (f_n(x) - f_m(x)) - (f_n(t) - f_m(t)) \right| &= |f'_n(t) - f'_m(t)| |x_0 - x| \\ &< \frac{\varepsilon}{2(b-a)} |x_0 - x| \\ &\leq \frac{\varepsilon}{2} \end{aligned} \tag{1}$$

Finally, by the triangle inequality,

$$\begin{aligned} |f_n(x) - f_m(x)| &\leq |f_n(x) - f_m(x) - f_n(x_0) + f_m(x_0)| + |f_n(x_0) - f_m(x_0)| \\ &< \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon. \end{aligned}$$

This holds true for all $x \in [a, b]$. Hence by 16.8, (f_n) converges uniformly on $[a, b]$.

Let

$$f(x) = \lim_{n \rightarrow \infty} f_n(x) \quad (a \leq x \leq b).$$

Fix $x \in [a, b]$, and let

$$\phi_n(t) = \frac{f_n(t) - f_n(x)}{t - x}, \quad \phi(t) = \frac{f(t) - f(x)}{t - x} \quad (a \leq t \leq b, t \neq x).$$

To show that f is differentiable, we need to show that $\lim_{t \rightarrow x} \phi(t)$ exists. Note that since f_n are differentiable, we have

$$\lim_{t \rightarrow x} \phi_n(t) = f'_n(x) \quad (n = 1, 2, \dots).$$

By (1), for all $n, m \geq N$,

$$|\phi_n(t) - \phi_m(t)| \leq \frac{\varepsilon}{2(b-a)},$$

so (ϕ_n) converges uniformly, for $t \neq x$. Since (f_n) converges to f , we conclude from (31) that

$$\lim_{n \rightarrow \infty} \phi_n(t) = \phi(t)$$

uniformly for $a \leq t \leq b, t \neq x$.

If we now apply Theorem 7.11 to (ϕ_n) , (32) and (33) show that

$$\lim_{t \rightarrow x} \phi(t) = \lim_{n \rightarrow \infty} f'_n(x),$$

and this is (27), by the definition of $\phi(t)$. □

Example 16.22 (Weierstrass function). Let us construct a continuous nowhere differentiable function on \mathbb{R} .

Define

$$\phi(x) = |x| \quad (-1 \leq x \leq 1).$$

We extend the definition of $\phi(x)$ to all of \mathbb{R} by making ϕ 2-periodic: $\phi(x) = \phi(x + 2)$. Then $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is continuous as $|\phi(x) - \phi(y)| \leq |x - y|$ (not hard to prove).

Let the *Weierstrass function* be defined as

$$f(x) = \sum_{n=0}^{\infty} \left(\frac{3}{4}\right)^n \phi(4^n x).$$

Claim. The Weierstrass function is continuous and nowhere differentiable on \mathbb{R} .

- Since $\sum \left(\frac{3}{4}\right)^n$ converges, and $|\phi(x)| \leq 1$ for all $x \in \mathbb{R}$, by the Weierstrass M-test, $f(x)$ converges uniformly and hence is continuous.
- Fix $x \in \mathbb{R}$ and $m \in \mathbb{Z}^+$, and define

$$\delta_m = \pm \frac{1}{2} \cdot 4^{-m},$$

where the sign is chosen in such a way so that there is no integer between $4^m x$ and $4^m(x + \delta_m)$, which can be done since $4^m |\delta_m| = \frac{1}{2}$. Define

$$\gamma_n = \frac{\phi(4^n(x + \delta_m)) - \phi(4^n x)}{\delta_m}.$$

If $n > m$, then as $4^n \delta_m$ is an even integer. Then as ϕ is 2-periodic we get that $\gamma_n = 0$.

Furthermore, since there is no integer between $4^m x \pm \frac{1}{2}$ and $4^m x$, we have that

$$\left| \phi\left(4^m x \pm \frac{1}{2}\right) - \phi(4^m x) \right| = \left| \left(4^m x \pm \frac{1}{2}\right) - 4^m x \right| = \frac{1}{2}.$$

Therefore

$$|\gamma_n| = \left| \frac{\phi\left(4^m x \pm \frac{1}{2}\right) - \phi(4^m x)}{\pm \frac{1}{2} \cdot 4^{-m}} \right| = 4^m.$$

Similarly, if $n < m$, since $|\phi(s) - \phi(t)| \leq |s - t|$,

$$|\gamma_n| = \left| \frac{\phi\left(4^n x \pm \frac{1}{2} \cdot 4^{n-m}\right) - \phi(4^n x)}{\pm \frac{1}{2} \cdot 4^{-m}} \right| \leq \left| \frac{\pm \frac{1}{2} \cdot 4^{n-m}}{\pm \frac{1}{2} \cdot 4^{-m}} \right| = 4^n.$$

Finally,

$$\begin{aligned} \left| \frac{f(x + \delta_m) - f(x)}{\delta_m} \right| &= \left| \sum_{n=0}^{\infty} \left(\frac{3}{4}\right)^n \frac{\phi(4^n(x + \delta_m)) - \phi(4^n x)}{\delta_m} \right| = \left| \sum_{n=0}^{\infty} \left(\frac{3}{4}\right)^n \gamma_n \right| \\ &= \left| \sum_{n=0}^m \left(\frac{3}{4}\right)^n \gamma_n \right| \\ &\geq \left| \frac{3^m}{4} \gamma_m \right| - \left| \sum_{n=0}^{m-1} \left(\frac{3}{4}\right)^n \gamma_n \right| \\ &\geq 3^m - \sum_{n=0}^{m-1} 3^n = 3^m - \frac{3^m - 1}{3 - 1} = \frac{3^m + 1}{2}. \end{aligned}$$

It is obvious that $\delta_m \rightarrow 0$ as $m \rightarrow \infty$, but $\frac{3^m + 1}{2}$ goes to infinity. Hence f cannot be differentiable at x .

§16.4 Equicontinuous Families of Functions

We would like an analogue of Bolzano–Weierstrass; that is, every bounded sequence of functions has a convergent subsequence.

Definition 16.23. Suppose (f_n) is a sequence of functions. We say (f_n) is *pointwise bounded* on E if for every $x \in E$, the sequence $(f_n(x))$ is bounded; that is,

$$\forall x \in E, \quad \exists M \in \mathbb{R}, \quad \forall n \in \mathbb{N}, \quad |f_n(x)| \leq M.$$

We say (f_n) is *uniformly bounded* on E if

$$\exists M \in \mathbb{R}, \quad \forall x \in E, n \in \mathbb{N}, \quad |f_n(x)| \leq M.$$

Lemma 16.24. Suppose (f_n) is a pointwise bounded sequence of complex-valued functions on a countable set E . Then (f_n) has a subsequence (f_{n_k}) such that $f_{n_k}(x)$ converges for every $x \in E$.

Proof. We will use a very common and useful diagonal argument.

Arrange the points of E in a sequence (x_i) , where $i = 1, 2, \dots$

Since (f_n) is pointwise bounded on E , the sequence $(f_n(x_1))_{n=1}^{\infty}$ is bounded. By the Bolzano–Weierstrass theorem, there exists a subsequence, which we denote by $(f_{1,k})_{k=1}^{\infty}$, such that $(f_{1,k}(x_1))_{k=1}^{\infty}$ converges.

Consider the array formed by the sequences S_1, S_2, \dots :

$$\begin{array}{cccc} S_1 : & f_{1,1} & f_{1,2} & f_{1,3} & \cdots \\ S_2 : & f_{2,1} & f_{2,2} & f_{2,3} & \cdots \\ S_3 : & f_{3,1} & f_{3,2} & f_{3,3} & \cdots \\ & \vdots & & & \end{array}$$

and which have the following properties:

- (i) S_n is a subsequence of S_{n-1} , for $n = 2, 3, \dots$
- (ii) $(f_{n,k}(x_n))$ converges, as $k \rightarrow \infty$ (the boundedness of $(f_n(x_n))$ makes it possible to choose S_n in this way);
- (iii) The order in which the functions appear is the same in each sequence; i.e., if one function precedes another in S_1 , they are in the same relation in every S_n , until one or the other is deleted. Hence, when going from one row in the above array to the next below, functions may move to the left but never to the right.

We now go down the diagonal of the array; i.e., we consider the sequence

$$S : f_{1,1} \quad f_{2,2} \quad f_{3,3} \quad \cdots$$

By (iii), the sequence S (except possibly its first $n - 1$ terms) is a subsequence of S_n , for $n = 1, 2, \dots$. Hence (ii) implies that $(f_{n,n}(x_i))$ converges, as $n \rightarrow \infty$, for every $x_i \in E$. □

to do

Definition 16.25. A family \mathcal{F} of functions $f : E \subset X \rightarrow \mathbb{C}$ is *equicontinuous* on E if

$$\forall \varepsilon > 0, \quad \exists \delta > 0, \quad \forall x, y \in E, f \in \mathcal{F}, \quad d(x, y) < \delta \implies |f(x) - f(y)| < \varepsilon.$$

Proposition 16.26. Suppose X is a compact metric space, $f_n \in \mathcal{C}(X, \mathbb{C})$, and (f_n) converges uniformly on X . Then (f_n) is equicontinuous on X .

Proof. Let $\varepsilon > 0$ be given. Since (f_n) converges uniformly on X , $f_n \rightarrow f$ on X with respect to the metric of $\mathcal{C}(X, \mathbb{C})$. Then

$$\lim_{n \rightarrow \infty} \|f_n - f\| = 0,$$

i.e., there exists $N \in \mathbb{N}$ such that for all $n \geq N$,

$$\|f_n - f_N\| < \frac{\varepsilon}{3}.$$

Since continuous functions are uniformly continuous on compact sets, f_n are uniformly continuous on K , so there exists $\delta > 0$ such that

$$d(x, y) < \delta \implies |f_i(x) - f_i(y)| < \frac{\varepsilon}{3}$$

for $i = 1, \dots, N$. If $n \geq N$ and $d(x, y) < \delta$,

$$\begin{aligned} |f_n(x) - f_n(y)| &\leq |f_n(x) - f_N(x)| + |f_N(x) - f_N(y)| + |f_N(y) - f_n(y)| \\ &< \frac{\varepsilon}{3} + \frac{\varepsilon}{3} + \frac{\varepsilon}{3} = \varepsilon. \end{aligned}$$

In conjunction with (43), this proves the theorem. □

We first need the following lemma.

Lemma 16.27. A compact metric space X contains a countable dense subset.

Proof. For each $n \in \mathbb{N}$, there exist finitely many balls of radius $\frac{1}{n}$ that cover X (by compactness of X). That is, for every n , there exist finitely many points $x_{n,1}, \dots, x_{n,k_n}$ such that

$$X = \bigcup_{i=1}^{k_n} B_{\frac{1}{n}}(x_{n,i}).$$

Claim. $S = \{x_{n,i} \mid i = 1, \dots, k_n\}$ is a countable dense subset of X .

- Since S is a countable union of finite sets, S is countable.
- For every $x \in X$ and every $\varepsilon > 0$, there exists $n \in \mathbb{N}$ such that $\frac{1}{n} < \varepsilon$ and an $x_{n,i} \in S$ such that

$$x \in B_{\frac{1}{n}}(x_{n,i}) \subset B_{\varepsilon}(x_{n,i}).$$

Hence $x \in \overline{S}$, so $\overline{S} = X$ and therefore S is dense. □

We can now prove the very useful Arzelà–Ascoli theorem about existence of convergent subsequences.

Theorem 16.28 (Arzelà–Ascoli theorem). *Suppose X is compact, $f_n \in \mathcal{C}(X, \mathbb{C})$, and (f_n) is pointwise bounded and equicontinuous on X . Then (f_n) is uniformly bounded on X , and contains a uniformly convergent subsequence.*

Proof. Let us first show that the sequence is uniformly bounded. By equicontinuity, there exists $\delta > 0$ such that

$$B_\delta(x) \subset f_n^{-1}(B_1(f_n(x))) \quad (x \in X).$$

Since X is compact, there exist finitely many points x_1, \dots, x_k such that

$$X = \bigcup_{j=1}^k B_\delta(x_j).$$

Since (f_n) is pointwise bounded, there exist M_1, \dots, M_k such that

$$|f_n(x_j)| \leq M_j \quad (j = 1, \dots, k)$$

for all n . Let $M = 1 + \max\{M_1, \dots, M_k\}$. Now given any $x \in X$, $x \in B_\delta(x_j)$ for some $1 \leq j \leq k$. Therefore, for all n we have $x \in f_n^{-1}(B_1(f_n(x_j)))$ or in other words

$$|f_n(x) - f_n(x_j)| < 1.$$

By reverse triangle inequality,

$$|f_n(x)| < 1 + |f_n(x_j)| \leq 1 + M_j \leq M$$

Since x was arbitrary, (f_n) is uniformly bounded.

Next, pick a countable dense set S . By Theorem 7.23, there exists a subsequence (f_{n_j}) that converges pointwise on S . Write $g_j = f_{n_j}$ for simplicity. Note that (g_n) is equicontinuous.

Let $\varepsilon > 0$ be given, then pick $\delta > 0$ such that for all $x \in X$,

$$B_\delta(x) \subset g_n^{-1}\left(B_{\frac{\varepsilon}{3}}(g_n(x))\right).$$

By density of S , every $x \in X$ is in some $B_\delta(y)$ for some $y \in S$, and by compactness of X , there is a finite subset $\{x_1, \dots, x_k\}$ of S such that

$$X = \bigcup_{j=1}^k B_\delta(x_j).$$

Now as there are finitely many points and we know that (g_n) converges pointwise on S , there exists $N \in \mathbb{N}$ such that for all $n, m \geq N$,

$$|g_n(x_j) - g_m(x_j)| < \frac{\varepsilon}{3} \quad (j = 1, \dots, k).$$

Let $x \in X$ be arbitrary. There is some i such that $x \in B_\delta(x_i)$ and so we have for all $i \in \mathbb{N}$,

$$|g_i(x) - g_i(x_j)| < \frac{\varepsilon}{3}$$

and so $n, m \geq N$ that

$$\begin{aligned} |g_n(x) - g_m(x)| &\leq |g_n(x) - g_n(x_j)| + |g_n(x_j) - g_m(x_j)| + |g_m(x_j) - g_m(x)| \\ &< \frac{\varepsilon}{3} + \frac{\varepsilon}{3} + \frac{\varepsilon}{3} = \varepsilon. \end{aligned}$$

□

Corollary 16.29. *Suppose X is a compact metric space. Let $S \subset \mathcal{C}(X, \mathbb{C})$ be a closed, bounded and equicontinuous set. Then S is compact.*

Corollary 16.30. *Suppose (f_n) is a sequence of differentiable functions on $[a, b]$, (f'_n) is uniformly bounded, and there exists $x_0 \in [a, b]$ such that $(f_n(x_0))$ is bounded. Then there exists a uniformly convergent subsequence (f_{n_k}) .*

§16.5 Stone–Weierstrass Approximation Theorem

Perhaps surprisingly, even a very badly behaving continuous function is really just a uniform limit of polynomials. We cannot really get any “nicer” as a function than a polynomial.

Weierstrass’s Version

Theorem 16.31 (Weierstrass approximation theorem). *If $f :: [a, b] \rightarrow \mathbb{C}$ is continuous, there exists a sequence of polynomials (P_n) such that $P_n \rightrightarrows f$ on $[a, b]$. If f is real, then P_n may be taken real.*

Proof. WLOG assume that $[a, b] = [0, 1]$. We may also assume that $f(0) = f(1) = 0$. For if the theorem is proved for this case, consider

$$g(x) = f(x) - f(0) - x[f(1) - f(0)] \quad (0 \leq x \leq 1).$$

Here $g(0) = g(1) = 0$, and if g can be obtained as the limit of a uniformly convergent sequence of polynomials, it is clear that the same is true for f , since $f - g$ is a polynomial.

Furthermore, we define $f(x)$ to be zero for x outside $[0, 1]$. Then f is uniformly continuous on the whole line.

Let

$$Q_n(x) = c_n(1 - x^2)^n \quad (n = 1, 2, \dots),$$

where c_n is chosen such that

$$\int_{-1}^1 Q_n(x) dx = 1 \quad (n = 1, 2, \dots).$$

We need some information about the order of magnitude of c_n . Since

$$\begin{aligned} \int_{-1}^1 (1 - x^2)^n dx &= 2 \int_0^1 (1 - x^2)^n dx \\ &\geq 2 \int_0^{\frac{1}{\sqrt{n}}} (1 - x^2)^n dx \\ &\geq 2 \int_0^{\frac{1}{\sqrt{n}}} (1 - nx^2) dx \\ &= \frac{4}{3\sqrt{n}} \\ &> \frac{1}{\sqrt{n}}, \end{aligned}$$

it follows from (48) that

$$c_n < \sqrt{n}.$$

The inequality $(1 - x^2)^n \geq 1 - nx^2$ which we used above is easily shown to be true by considering the function

$$(1 - x^2)^n - 1 + nx^2$$

which is zero at $x = 0$ and whose derivative is positive in $(0, 1)$.

For any $\delta > 0$, (49) implies

$$Q_n(x) \leq \sqrt{n}(1 - \delta^2)^n \quad (\delta \leq |x| \leq 1),$$

so that $Q_n \Rightarrow 0$ in $\delta \leq |x| \leq 1$.

Now let

$$P_n(x) = \int_{-1}^1 f(x+t)Q_n(t) dt \quad (0 \leq x \leq 1).$$

Our assumptions about f show, by a simple change of variable, that

$$P_n(x) = \int_{-x}^{1-x} f(x+t)Q_n(t) dt = \int_0^1 f(t)Q_n(t-x) dt,$$

and the last integral is clearly a polynomial in x . Thus (P_n) is a sequence of polynomials, which are real if f is real.

Given $\varepsilon > 0$, we choose $\delta > 0$ such that

$$|y-x| < \delta \implies |f(y) - f(x)| < \frac{\varepsilon}{2}.$$

Let $M = \sup |f(x)|$. Using (48), (50), and the fact that $Q_n(x) \geq 0$, we see that for $0 \leq x \leq 1$,

$$\begin{aligned} |P_n(x) - f(x)| &= \left| \int_{-1}^1 [f(x+t) - f(x)]Q_n(t) dt \right| \\ &\leq \int_{-1}^1 |f(x+t) - f(x)|Q_n(t) dt \\ &\leq 2M \int_{-1}^{-\delta} Q_n(t) dt + \frac{\varepsilon}{2} \int_{-\delta}^{\delta} Q_n(t) dt + 2M \int_{\delta}^1 Q_n(t) dt \\ &\leq 4M\sqrt{n}(1 - \delta^2)^n + \frac{\varepsilon}{2} \\ &< \varepsilon \end{aligned}$$

for all large enough n , which proves the theorem. □

Think about the consequences of the theorem. If you have any property that gets preserved under uniform convergence and it is true for polynomials, then it must be true for all continuous functions.

Let us note an immediate application of the Weierstrass theorem. We have already seen that countable dense subsets can be very useful.

Corollary 16.32. *The metric space $\mathcal{C}([a, b], \mathbb{C})$ contains a countable dense subset.*

Corollary 16.33. *For every interval $[-a, a]$, there exists a sequence of real polynomials P_n such that $P_n(0) = 0$ and*

$$\lim_{n \rightarrow \infty} P_n(x) = |x|$$

uniformly on $[-a, a]$.

Algebra of Functions

We shall now isolate those properties of the polynomials which make the Weierstrass theorem possible.

Definition 16.34. A family \mathcal{A} of complex-valued functions $f : X \rightarrow \mathbb{C}$ is an *algebra* if, for all $f, g \in \mathcal{A}, c \in \mathbb{C}$,

(i) $f + g \in \mathcal{A}$; (closed under addition)

(ii) $fg \in \mathcal{A}$; (closed under multiplication)

(iii) $cf \in \mathcal{A}$. (closed under scalar multiplication)

If we talk of an algebra of real-valued functions, then of course we only need the above to hold for $c \in \mathbb{R}$.

\mathcal{A} is *uniformly closed* if the limit of every uniformly convergent sequence in \mathcal{A} is also in \mathcal{A} .

Let \mathcal{B} be the set of all limits of uniformly convergent sequences in \mathcal{A} . Then \mathcal{B} is the *uniform closure* of \mathcal{A} .

Example 16.35.

- $\mathcal{C}(X, Y)$ is an algebra of functions.

Proposition 16.36. Let \mathcal{B} be the uniform closure of an algebra \mathcal{A} of bounded functions. Then \mathcal{B} is a uniformly closed algebra.

Now let us distill the right properties of polynomials that were sufficient for an approximation theorem.

Definition 16.37. Let \mathcal{A} be a family of functions defined on X .

We say \mathcal{A} *separates points* if for every $x, y \in X$, with $x \neq y$ there exists $f \in \mathcal{A}$ such that $f(x) \neq f(y)$.

We say \mathcal{A} *vanishes at no point* if for every $x \in X$ there exists $f \in \mathcal{A}$ such that $f(x) \neq 0$.

Example 16.38.

Proposition 16.39. Suppose \mathcal{A} is an algebra of functions on X , that separates points and vanishes at no point. Suppose x, y are distinct points of X and $c, d \in \mathbb{C}$. Then there exists $f \in \mathcal{A}$ such that

$$f(x) = c, \quad f(y) = d.$$

The Theorem

We now have all the material needed for Stone's generalisation of the Weierstrass theorem.

Theorem 16.40 (Stone–Weierstrass approximation theorem). Let X be a compact metric space and \mathcal{A} an algebra of real-valued continuous functions on X , such that \mathcal{A} separates points and vanishes at no point. Then the uniform closure of \mathcal{A} is all of $\mathcal{C}(X, \mathbb{R})$.

Exercises

17 Some Special Functions

§17.1 Power Series

Definition 17.1. Given a sequence (c_n) of complex numbers, a *power series* takes the form

$$\sum_{n=0}^{\infty} c_n z^n,$$

where $z \in \mathbb{C}$; the numbers c_n are called the *coefficients* of the series.

The convergence of $\sum c_n z^n$ depends on the choice of z (we would expect that a power series will be more likely to converge for small $|z|$ than for large $|z|$). More specifically, there is a “circle of convergence”, where $\sum c_n z^n$ converges if z is in the interior of the circle, and diverges if z is in the exterior.

Lemma 17.2 (Cauchy–Hadamard theorem). *Given the power series $\sum c_n z^n$, let*

$$\alpha = \limsup_{n \rightarrow \infty} \sqrt[n]{|c_n|}, \quad R = \frac{1}{\alpha}.$$

(If $\alpha = 0$, $R = +\infty$; if $\alpha = +\infty$, $R = 0$.) Then $\sum c_n z^n$

(i) converges if $|z| < R$,

(ii) diverges if $|z| > R$.

R is called the *radius of convergence* of $\sum c_n (z - a)^n$; the *disk of convergence* for the power series is

$$D_R(a) := \{z \in \mathbb{C} : |z - a| < R\}.$$

Proof. Let $a_n = c_n z^n$. We apply the root test:

$$\limsup_{n \rightarrow \infty} \sqrt[n]{|a_n|} = \limsup_{n \rightarrow \infty} \sqrt[n]{|c_n z^n|} = |z| \limsup_{n \rightarrow \infty} \sqrt[n]{|c_n|} = \frac{|z|}{R}.$$

(i) If $|z| < R$, then $\limsup_{n \rightarrow \infty} \sqrt[n]{|a_n|} < 1$. By the root test, $\sum c_n z^n$ converges absolutely and thus converges.

(ii) If $|z| > R$, then $\limsup_{n \rightarrow \infty} \sqrt[n]{|a_n|} > 1$. By the root test, $\sum c_n z^n$ diverges.

□

Example 17.3.

- $\sum_{k=0}^{\infty} z^k$ has radius of convergence 1.

- $\sum_{k=0}^{\infty} \frac{1}{n^n} z^k$ has radius of convergence ∞ .
- $\sum_{k=0}^{\infty} \frac{1}{n!} z^k$ has radius of convergence ∞ .
- $\sum_{k=0}^{\infty} n^n z^k$ has radius of convergence 0, so it converges only if $z = 0$.

In the previous result, we have shown that the radius of convergence can be found by using the root test. We can also find it using the ratio test (which is easier to compute).

Lemma 17.4. *If $\sum c_n z^n$ has radius of convergence R , then*

$$R = \lim_{n \rightarrow \infty} \left| \frac{c_n}{c_{n+1}} \right|,$$

if this limit exists.

Proof. By the ratio test, $\sum c_n z^n$ converges if

$$\lim_{n \rightarrow \infty} \left| \frac{c_{n+1} z^{n+1}}{c_n z^n} \right| < 1.$$

This is equivalent to

$$|z| < \frac{1}{\lim_{n \rightarrow \infty} \left| \frac{c_{n+1}}{c_n} \right|} = \lim_{n \rightarrow \infty} \left| \frac{c_n}{c_{n+1}} \right|.$$

□

Proposition 17.5. *Suppose the radius of convergence of $\sum c_n z^n$ is 1, and suppose $c_0 \geq c_1 \geq c_2 \geq \dots$, $c_n \rightarrow 0$. Then $\sum c_n z^n$ converges at every point on the circle $|z| = 1$, except possibly at $z = 1$.*

Proof. Let

$$a_n = z^n, \quad b_n = c_n.$$

Then the hypothesis of Proposition 12.44 are satisfied, since

$$|A_n| = \left| \sum_{k=0}^n z^k \right| = \left| \frac{1 - z^{n+1}}{1 - z} \right| \leq \frac{2}{|1 - z|}$$

if $|z| = 1$, $|z| \neq 1$.

□

Definition 17.6. An *analytic function* is a function that can be represented by a power series; that is, functions of the form

$$f(x) = \sum_{n=0}^{\infty} c_n x^n$$

or, more generally,

$$f(x) = \sum_{n=0}^{\infty} c_n(x-a)^n.$$

We shall restrict ourselves to real values of x (since we have yet to define complex differentiation). Instead of circles of convergence we shall therefore encounter intervals of convergence.

As a matter of convenience, we shall often take $a = 0$ without any loss of generality. If $\sum c_n x^n$ converges for all $x \in (-R, R)$, for some $R > 0$, we say that f is *expanded in a power series* about the point $x = 0$.

Proposition 17.7. *Suppose $\sum c_n x^n$ converges for $|x| < R$. Let*

$$f(x) = \sum_{n=0}^{\infty} c_n x^n \quad (|x| < R).$$

Then

(i) $\sum c_n x^n$ converges absolutely and uniformly on $(-r, r)$ where $r < R$;

(ii) $f(x)$ is continuous and differentiable on $(-R, R)$, and

$$f'(x) = \sum_{n=1}^{\infty} n c_n x^{n-1} \quad (|x| < R).$$

Proof.

(i) We will show that $\sum c_n x^n$ converges absolutely and uniformly on $[-R + \varepsilon, R - \varepsilon]$ for all $\varepsilon > 0$.

Idea. Weierstrass M-test.

Let $\varepsilon > 0$ be given. For $|x| \leq R - \varepsilon$, notice that we have

$$|c_n x^n| \leq |c_n|(R - \varepsilon)^n \quad (n = 1, 2, \dots).$$

Consider the series

$$\sum |c_n|(R - \varepsilon)^n.$$

By Lemma 17.2, which states that every power series converges (absolutely) in the interior of its interval of convergence, we have that $\sum |c_n|(R - \varepsilon)^n$ converges.

By the Weierstrass M-test, $\sum c_n x^n$ uniformly converges on $[-R + \varepsilon, R - \varepsilon]$.

(ii) Since $\lim_{n \rightarrow \infty} \sqrt[n]{n} = 1$, we have

$$\limsup_{n \rightarrow \infty} \sqrt[n]{n|c_n|} = \limsup_{n \rightarrow \infty} \sqrt[n]{|c_n|},$$

so the series $\sum_{n=0}^{\infty} c_n x^n$ and $\sum_{n=1}^{\infty} n c_n x^{n-1}$ have the same radius of convergence; thus $\sum_{n=1}^{\infty} n c_n x^{n-1}$ has radius of convergence R .

Idea. Interchange limits and derivatives.

Since $\sum_{n=1}^{\infty} nc_n x^{n-1}$ is a power series, by (i), it converges uniformly in $[-R + \varepsilon, R - \varepsilon]$, for every $\varepsilon > 0$.

Consider the partial sums $f_n(x) = \sum_{k=0}^n c_k x^k$; evidently f_n are differentiable on $[-R + \varepsilon, R - \varepsilon]$, and $(f_n(x))$ converges whenever $|x| < R$. Also,

$$f'_n(x) = \sum_{k=1}^n kc_k x^{k-1},$$

converge uniformly on $[-R + \varepsilon, R - \varepsilon]$. By Proposition 16.21, for all $|x| < R$,

$$f'(x) = \lim_{n \rightarrow \infty} f'_n(x) = \lim_{n \rightarrow \infty} \sum_{k=1}^n kc_k x^{k-1} = \sum_{n=1}^{\infty} nc_n x^{n-1}.$$

Since f is differentiable on $(-R, R)$, by Lemma 14.2, f is continuous on $(-R, R)$.

□

Corollary 17.8. *f is infinitely differentiable in $(-R, R)$; its derivatives are given by*

$$f^{(k)}(x) = \sum_{n=k}^{\infty} n(n-1)\cdots(n-k+1)c_n x^{n-k}. \tag{17.1}$$

In particular,

$$f^{(k)}(0) = k!c_k \quad (k = 0, 1, 2, \dots). \tag{17.2}$$

Proof. Apply the previous result successively to f, f', f'', \dots

Then plug in $x = 0$.

□

Remark. Eq. (17.2) is very interesting.

- It shows, on the one hand, that the coefficients of the power series development of f are determined by the values of f and its derivatives at a single point.
- On the other hand, if the coefficients are given, the values of the derivatives of f at the center of the interval of convergence can be read off immediately from the power series.

If the series (3) converges at an endpoint, say at $x = R$, then f is continuous not only in $(-R, R)$, but also at $x = R$, as shown by the following result (for simplicity of notation, we take $R = 1$).

Proposition 17.9 (Abel's theorem). *Suppose $\sum c_n$ converges. Let*

$$f(x) = \sum_{n=0}^{\infty} c_n x^n \quad (-1 < x < 1).$$

Then $f(x)$ is continuous at $x = 1$.

Proof. We want to show that $\lim_{x \rightarrow 1} f(x) = \sum_{n=0}^{\infty} c_n$.

Let

$$s_n = c_0 + \cdots + c_n, \quad s_{-1} = 0.$$

Then we write

$$\begin{aligned} \sum_{n=0}^m c_n x^n &= \sum_{n=0}^m (s_n - s_{n-1}) x^n \\ &= \sum_{n=0}^m s_n x^n - \sum_{n=1}^m s_{n-1} x^n \\ &= \sum_{n=0}^m s_n x^n - \sum_{n=0}^{m-1} s_n x^{n+1} \\ &= (1-x) \sum_{n=0}^{m-1} s_n x^n + s_m x^m. \end{aligned}$$

For $|x| < 1$, we let $m \rightarrow \infty$ and obtain

$$f(x) = (1-x) \sum_{n=0}^{\infty} s_n x^n.$$

Suppose $s_n \rightarrow s$. We will show that $\lim_{x \rightarrow 1} f(x) = s$. Fix $\varepsilon > 0$, there exists $N \in \mathbb{N}$ such that

$$n \geq N \implies |s_n - s| < \frac{\varepsilon}{2}.$$

Note that for $|x| < 1$, since $\sum_{n=0}^{\infty} x^n = \frac{1}{1-x}$, we can write

$$(1-x) \sum_{n=0}^{\infty} s x^n = s.$$

If $x > 1 - \delta$, for some suitably chosen $\delta > 0$, we have

$$\begin{aligned} |f(x) - s| &= \left| (1-x) \sum_{n=0}^{\infty} (s_n - s) x^n \right| \\ &= (1-x) \left| \sum_{n=0}^N (s_n - s) x^n + \sum_{n=N+1}^{\infty} (s_n - s) x^n \right| \\ &\leq (1-x) \left| \sum_{n=0}^N (s_n - s) x^n \right| + (1-x) \sum_{n=N+1}^{\infty} (s_n - s) x^n. \end{aligned}$$

Note that

$$\begin{aligned} (1-x) \left| \sum_{n=N+1}^{\infty} (s_n - s)x^n \right| &\leq (1-x) \sum_{n=N+1}^{\infty} |s_n - s| |x|^n \\ &< \frac{\varepsilon}{2} (1-x) \sum_{n=N+1}^{\infty} x^n \\ &= \frac{\varepsilon}{2} (1-x) \frac{x^{N+1}}{1-x} < \frac{\varepsilon}{2} \end{aligned}$$

and

$$\begin{aligned} (1-x) \left| \sum_{n=0}^N (s_n - s)x^n \right| &\leq (1-x) \sum_{n=0}^N |s_n - s| |x|^n \\ &< (1-x) \sum_{n=0}^N |s_n - s| \end{aligned}$$

which can be bounded by, say, M because there are only finitely many terms in the sum. Choosing $\delta < \frac{\varepsilon}{2M}$ gives

$$(1-x) \sum_{n=0}^N |s_n - s| < (1 - (1 - \delta)) \sum_{n=0}^N |s_n - s| < \delta M < \frac{\varepsilon}{2}.$$

Hence

$$|f(x) - s| < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon,$$

and therefore $\lim_{x \rightarrow 1} f(x) = s$, as desired. □

We now require a result concerning an inversion in the order of summation.

Proposition 17.10 (Fubini's theorem for sums). *Given a double sequence (a_{ij}) , $i = 1, 2, \dots$, $j = 1, 2, \dots$, suppose that*

$$\sum_{j=1}^{\infty} |a_{ij}| = b_i \quad (i = 1, 2, \dots)$$

and $\sum b_i$ converges. Then

$$\sum_{i=1}^{\infty} \sum_{j=1}^{\infty} a_{ij} = \sum_{j=1}^{\infty} \sum_{i=1}^{\infty} a_{ij}. \tag{17.3}$$

Remark. This is analogous to Fubini's theorem for the swapping of double integrals.

Proof. Let E be a countable set:

$$E = \{x_0, x_1, x_2, \dots\},$$

and suppose $x_n \rightarrow x_0$ as $n \rightarrow \infty$. Define the sequence of functions $f_i : E \rightarrow \mathbb{C}$ by

$$\begin{aligned} f_i(x_0) &= \sum_{j=1}^{\infty} a_{ij} \quad (i = 1, 2, \dots) \\ f_i(x_n) &= \sum_{j=1}^n a_{ij} \quad (i, n = 1, 2, \dots) \end{aligned}$$

See that each f_i is continuous at x_0 .

Since $|f_i(x)| \leq b_i$ for $x \in E$ (by triangle inequality), and $\sum b_i$ converges, by the Weierstrass M-test, $\sum_{i=1}^n f_i(x)$ converges uniformly. Let

$$g(x) = \sum_{i=1}^{\infty} f_i(x) \quad (x \in E).$$

By 7.11, g is continuous at x_0 , so

$$\lim_{n \rightarrow \infty} g(x_n) = g(x_0).$$

It follows that

$$\begin{aligned} \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} a_{ij} &= \sum_{i=1}^{\infty} f_i(x_0) = g(x_0) = \lim_{n \rightarrow \infty} g(x_n) \\ &= \lim_{n \rightarrow \infty} \sum_{i=1}^{\infty} f_i(x_n) = \lim_{n \rightarrow \infty} \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} a_{ij} \\ &= \lim_{n \rightarrow \infty} \sum_{j=1}^n \sum_{i=1}^{\infty} a_{ij} = \sum_{j=1}^{\infty} \sum_{i=1}^{\infty} a_{ij}. \end{aligned}$$

□

Theorem 17.11 (Taylor's theorem). *Suppose $\sum c_n x^n$ converges in $|x| < R$, let*

$$f(x) = \sum_{n=0}^{\infty} c_n x^n.$$

If $a \in (-R, R)$, then f can be expanded in a power series about the point $x = a$ which converges in $|x - a| < R - |a|$, and

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x - a)^n \quad (|x - a| < R - |a|). \tag{17.4}$$

Proof. We have

$$\begin{aligned} f(x) &= \sum_{n=0}^{\infty} c_n x^n = \sum_{n=0}^{\infty} c_n [(x - a) + a]^n \\ &= \sum_{n=0}^{\infty} c_n \sum_{m=0}^n \binom{n}{m} a^{n-m} (x - a)^m \\ &= \sum_{m=0}^{\infty} \left[\sum_{n=m}^{\infty} \binom{n}{m} c_n a^{n-m} \right] (x - a)^m \end{aligned} \tag{1}$$

This is the desired expansion about the point $x = a$. We need to show that the swapping of summations in (1) is valid, which is applicable only if $\sum_{n=m}^{\infty} \binom{n}{m} c_n a^{n-m} (x - a)^m$ satisfies Theorem 8.3, i.e.

$$\sum_{n=0}^{\infty} \sum_{m=0}^n \left| c_n \binom{n}{m} a^{n-m} (x - a)^m \right|$$

converges. Write

$$\begin{aligned} & \sum_{n=0}^{\infty} \sum_{m=0}^n \left| c_n \binom{n}{m} a^{n-m} (x-a)^m \right| \\ &= \sum_{n=0}^{\infty} \sum_{m=0}^{\infty} \binom{n}{m} |c_n| |a|^{n-m} |x-a|^m \\ &= \sum_{n=0}^{\infty} |c_n| (|x-a| + |a|)^n, \end{aligned}$$

which converges if and only if $|x-a| + |a| < R$.

Finally, the form of the coefficients in Eq. (17.4) follows from Eq. (17.2): differentiate $f(x)$ repeatedly to obtain

$$\begin{aligned} f^{(m)}(x) &= \sum_{n=m}^{\infty} c_n n(n-1) \cdots (n-m+1) x^{n-m} \\ &= \sum_{n=m}^{\infty} c_n m! \binom{n}{m} x^{n-m} \end{aligned}$$

and then plug in $x = a$. □

If two power series converge to the same function in $(-R, R)$, (7) shows that the two series must be identical, i.e., they must have the same coefficients. It is interesting that the same conclusion can be deduced from much weaker hypotheses:

Proposition 17.12. Suppose $\sum a_n x^n$ and $\sum b_n x^n$ converge in $S = (-R, R)$. Let E be the set of all $x \in S$ such that

$$\sum_{n=0}^{\infty} a_n x^n = \sum_{n=0}^{\infty} b_n x^n.$$

If E has a limit point in S , then $a_n = b_n$ for $n = 0, 1, 2, \dots$. Hence (20) holds for all $x \in S$.

Proof. Let $c_n = a_n - b_n$, and let

$$f(x) = \sum_{n=0}^{\infty} c_n x^n \quad (x \in S)$$

We will show that $c_n = 0$, so that $f(x) = 0$ on E .

Let A be the set of all limit points of E in S , and let B consist of all other points of S . It is clear from the definition of "limit point" that B is open. Suppose we can prove that A is open. Then A and B are disjoint open sets. Hence they are separated (Definition 2.45). Since $S = A \cup B$, and S is connected, one of A and B must be empty. By hypothesis, A is not empty. Hence B is empty, and $A = S$. Since f is continuous in S , $A \subset E$. Thus $E = S$, and (7) shows that $c_n = 0$ for $n = 0, 1, 2, \dots$, which is the desired conclusion.

Thus we have to prove that A is open. If $x_0 \in A$, Theorem 8.4 shows that

$$f(x) = \sum_{n=0}^{\infty} d_n (x - x_0)^n \quad (|x - x_0| < R - |x_0|).$$

We claim that $d_n = 0$ for all n . Otherwise, let k be the smallest nonnegative integer such that $d_k \neq 0$. Then

$$f(x) = (x - x_0)^k g(x) \quad (|x - x_0| < R - |x_0|),$$

where

$$g(x) = \sum_{m=0}^{\infty} d_{k+m} (x - x_0)^m.$$

Since g is continuous at x_0 and

$$g(x_0) = d_k \neq 0,$$

there exists $\delta > 0$ such that $g(x) \neq 0$ if $|x - x_0| < \delta$. It follows from (23) that $f(x) \neq 0$ if $0 < |x - x_0| < \delta$. But this contradicts the fact that x_0 is a limit point of E .

Thus $d_n = 0$ for all n , so that $f(x) = 0$ for all x for which (22) holds, i.e., in a neighborhood of x_0 . This shows that A is open, and completes the proof.

to do

□

Exponential and Logarithmic Functions

Definition 17.13 (Exponential function). For $z \in \mathbb{C}$, define

$$\exp(z) := \sum_{n=0}^{\infty} \frac{z^n}{n!}. \quad (17.5)$$

Proposition 17.14. $\exp(z)$ converges for every $z \in \mathbb{C}$.

Proof. Ratio test. □

The series (1) converges absolutely for every z and converges uniformly on every bounded subset of the complex plane. Thus \exp is a continuous function.

Lemma 17.15 (Addition formula). For $z, w \in \mathbb{C}$,

$$\exp(z + w) = \exp(z) \exp(w). \quad (17.6)$$

Proof. By multiplication of absolutely convergent series, we have

$$\begin{aligned} \exp(z) \exp(w) &= \sum_{n=0}^{\infty} \frac{z^n}{n!} \sum_{m=0}^{\infty} \frac{w^m}{m!} \\ &= \sum_{k=0}^{\infty} \left(\frac{z^k}{k!} + \frac{z^{k-1}}{(k-1)!} \frac{w}{1!} + \cdots + \frac{w^k}{k!} \right) \\ &= \sum_{k=0}^{\infty} \frac{1}{k!} \sum_{m+n=k} \binom{k}{n} z^n w^{k-n} \\ &= \sum_{k=0}^{\infty} \frac{1}{k!} (z + w)^k \\ &= \exp(z + w) \end{aligned}$$

□

Corollary 17.16. For $z \in \mathbb{C}$,

$$\exp(z) \exp(-z) = 1.$$

Proof. Take $z = z, w = -z$ in Eq. (17.6). □

Evidently $e = \exp(1)$, and we shall usually replace $\exp(z)$ by the customary shorter expression e^z . Note that $e^0 = \exp(0) = 1$, by (1).

Lemma 17.17 (Basic properties of \exp).

- (i) $\exp(z) \neq 0$ for every $z \in \mathbb{C}$.
- (ii) \exp is its own derivative: $\exp'(z) = \exp(z)$

(iii) The restriction of \exp to \mathbb{R} is a monotonically increasing positive function, and

$$\lim_{x \rightarrow \infty} e^x = \infty, \quad \lim_{x \rightarrow -\infty} e^x = 0.$$

(iv) There exists a positive number π such that $e^{\frac{\pi i}{2}} = i$ and such that $e^z = 1$ if and only if $\frac{z}{2\pi i}$ is an integer.

(v) \exp is a periodic function, with period $2\pi i$.

(vi) The mapping $t \mapsto e^{it}$ maps the real axis onto the unit circle.

(vii) If $w \in \mathbb{C}$, $w \neq 0$, then $w = e^z$ for some z .

Proof.

(i) This follows from the previous corollary.

(ii) We have

$$\exp'(z) = \lim_{h \rightarrow 0} \frac{\exp(z+h) - \exp(z)}{h} = \exp(z) \lim_{h \rightarrow 0} \frac{\exp(h) - 1}{h} = \exp(z),$$

where the first equality is a matter of definition, the second follows from Eq. (17.6), and the third from (1).

(iii)

□

We shall encounter the integral of $(1+x^2)^{-1}$ over the real line. To evaluate it, put $\phi(t) = \frac{\sin t}{\cos t}$ in $(-\frac{\pi}{2}, \frac{\pi}{2})$. By (6), $\phi' = 1 + \phi^2$. Hence ϕ is a monotonically increasing mapping of $(-\frac{\pi}{2}, \frac{\pi}{2})$ onto $(-\infty, \infty)$, and we obtain

$$\int_{-\infty}^{\infty} \frac{1}{1+x^2} dx = \int_{-\frac{\pi}{2}}^{\frac{\pi}{2}} \frac{\phi'(t)}{1+\phi^2(t)} dt = \int_{-\frac{\pi}{2}}^{\frac{\pi}{2}} dt = \pi.$$

Trigonometric Functions

Definition 17.18. For $z \in \mathbb{C}$, define

$$\cos z := \frac{e^{iz} + e^{-iz}}{2}, \quad \sin z := \frac{e^{iz} - e^{-iz}}{2i}. \quad (17.7)$$

By Eq. (17.5), we obtain the power series

$$\begin{aligned} \cos z &= \sum_{n=0}^{\infty} (-1)^n \frac{z^{2n}}{(2n)!} \\ \sin z &= \sum_{n=0}^{\infty} (-1)^n \frac{z^{2n+1}}{(2n+1)!} \end{aligned}$$

Lemma 17.19 (Euler's identity).

$$e^{iz} = \cos z + i \sin z.$$

Proof. This is immediate from Eq. (17.7). □

Thus when x is real, we note that from definition

$$\cos x = \operatorname{Re} e^{ix}, \quad \sin x = \operatorname{Im} e^{ix}.$$

In other words, sine and cosine are real-valued when we plug in real x .

Lemma 17.20 (Basic properties).

- (i) For $x \in \mathbb{R}$, $\cos' x = -\sin x$ and $\sin' x = \cos x$.
- (ii) \exp is periodic, with period $2\pi i$.
- (iii) C and S are periodic, with period 2π .
- (iv) If $0 < t < 2\pi$, then $\exp(it) \neq 1$.
- (v) If $z \in \mathbb{C}$, $|z| = 1$, there exists a unique $t \in [0, 2\pi)$ such that $\exp(it) = z$.

§17.2 Algebraic Completeness of the Complex Field

We now prove that the complex field is *algebraically complete*; that is, every non-constant polynomial with complex coefficients has a complex root.

Theorem 17.21 (Fundamental Theorem of Algebra). *For $a_i \in \mathbb{C}$, let*

$$P(z) = \sum_{k=0}^n a_k z^k$$

where $n \geq 1$, $a_n \neq 0$. Then $P(z) = 0$ for some $z \in \mathbb{C}$.

Proof. WLOG assume $a_n = 1$. Let

$$\mu = \inf |P(z)| \quad (z \in \mathbb{C}).$$

If $|z| = R$, then

$$|P(z)| \geq R^n \left(1 - |a_{n-1}|R^{-1} - \cdots - |a_0|R^{-n} \right).$$

The RHS tends to ∞ as $R \rightarrow \infty$. Hence there exists R_0 such that $|P(z)| > \mu$ if $|z| > R_0$. Since $|P|$ is continuous on the closed disk $\overline{D}_{R_0}(0)$, Theorem 4.16 shows that $|P(z_0)| = \mu$ for some z_0 .

Claim. $\mu = 0$.

If not, let $Q(z) = \frac{P(z + z_0)}{P(z_0)}$. Then Q is a non-constant polynomial, $Q(0) = 1$, and $|Q(z)| \geq 1$ for all z . There is a smallest integer k , $1 \leq k \leq n$ such that

$$Q(z) = 1 + b_k z^k + \cdots + b_n z^n \quad (b_k \neq 0).$$

By Theorem 8.7(d) there is a real θ such that

$$e^{ik\theta} b_k = -|b_k|.$$

If $r > 0$ and $r^k |b_k| < 1$, the above equation implies

$$|1 + b_k r^k e^{ik\theta}| = 1 - r^k |b_k|,$$

so that

$$\left| Q\left(r e^{i\theta} \right) \right| \leq 1 - r^k \left(|b_k| - r |b_{k+1}| - \cdots - r^{n-k} |b_n| \right).$$

For sufficiently small r , the expression in braces is positive; hence $|Q(r e^{i\theta})| < 1$, a contradiction.

Thus $\mu = 0$, that is, $P(z_0) = 0$. □

§17.3 Fourier Series

Definition 17.22. A *trigonometric polynomial* is a finite sum of the form

$$f(x) = a_0 + \sum_{n=1}^N (a_n \cos nx + b_n \sin nx) \quad (x \in \mathbb{R})$$

where $a_0, a_1, \dots, a_N, b_1, \dots, b_N \in \mathbb{C}$.

Using Eq. (17.7), we can write the above in the form

$$f(x) = \sum_{n=-N}^N c_n e^{inx}$$

for some constants $c_n \in \mathbb{C}$. This is a more convenient form of trigonometric polynomials, which we shall work with.

It is clear that every trigonometric polynomial is periodic, with period 2π .

For non-zero integer n , e^{inx} is the derivative of $\frac{1}{in}e^{inx}$, which also has period 2π . Hence

$$\frac{1}{2\pi} \int_{-\pi}^{\pi} e^{inx} dx = \begin{cases} 1 & (n = 0) \\ 0 & (n = \pm 1, \pm 2, \dots) \end{cases}$$

Definition 17.23. Let $f \in \mathcal{R}[-\pi, \pi]$. The *Fourier coefficients* of f are the numbers c_n , defined by

$$c_n = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) e^{-inx} dx.$$

The series

$$\sum_{n=-\infty}^{\infty} c_n e^{inx}$$

formed with the Fourier coefficients is called the *Fourier series* of f ; in this case we write

$$f \sim \sum_{n=-\infty}^{\infty} c_n e^{inx}.$$

We say f is an L^2 function if $|f|^2$ is Lebesgue integrable. The space of L^2 functions on a set E is denoted by $L^2(E)$. For all $f, g \in L^2(E)$, define the inner product

$$\langle f, g \rangle = \int_E f(x) \overline{g(x)} dx.$$

Then the norm of f squared is defined as

$$\|f\|^2 := \langle f, f \rangle = \int_E |f(x)|^2 dx.$$

We say that f and g are *orthogonal* if $\langle f, g \rangle = 0$.

Definition 17.24. Let (ϕ_n) be a sequence of complex functions on $[a, b]$.

- (i) We say (ϕ_n) is an **orthogonal system** of functions on $[a, b]$ if $\langle \phi_n, \phi_m \rangle = 0$ for all $n \neq m$.
- (ii) We say (ϕ_n) is an **orthonormal system** of functions on $[a, b]$ if (ϕ_n) is an orthogonal system, and $\|\phi_n\| = 1$ for all n .

Example 17.25.

- $\left\{ \frac{1}{\sqrt{2\pi}} e^{inx} \right\}$ is an orthonormal system on $[-\pi, \pi]$.
- $\left\{ \frac{1}{\sqrt{2\pi}}, \frac{1}{\sqrt{\pi}} \cos nx, \frac{1}{\sqrt{\pi}} \sin nx \right\}$ is an orthonormal system on $[-\pi, \pi]$.

Proof. We have

$$\begin{aligned} & \int_{-\pi}^{\pi} \cos nx \sin mx \, dx \\ &= \int_{-\pi}^{\pi} \frac{e^{inx} + e^{-inx}}{2} \frac{e^{imx} - e^{-imx}}{2i} \, dx \\ &= \int_{-\pi}^{\pi} \frac{e^{i(n+m)x} - e^{i(n-m)x} + e^{-i(n-m)x} - e^{-i(n+m)x}}{4i} \, dx = 0 \end{aligned}$$

and similarly

$$\begin{aligned} & \int_{-\pi}^{\pi} \cos nx \cos mx \, dx \\ &= \int_{-\pi}^{\pi} \frac{e^{inx} + e^{-inx}}{2} \frac{e^{imx} + e^{-imx}}{2i} \, dx \\ &= \int_{-\pi}^{\pi} \frac{e^{i(n+m)x} + e^{i(n-m)x} + e^{-i(n-m)x} + e^{-i(n+m)x}}{4} \, dx \\ &= \begin{cases} \frac{1}{2} \cdot 2\pi = \pi & (n = m) \\ 0 & (n \neq m) \end{cases} \end{aligned}$$

□

If (ϕ_n) is an orthonormal system of functions on $[a, b]$, then

$$f \sim \sum_{n=1}^{\infty} c_n \phi_n$$

where $c_n = \langle f, \phi_n \rangle$; we call c_n the n -th Fourier coefficient of f relative to (ϕ_n) .

Example 17.26. In \mathbb{R}^3 , let

$$\phi_1 = (1, 0, 0), \quad \phi_2 = (0, 1, 0), \quad \phi_3 = (0, 0, 1).$$

Suppose $f = (2, -1, 3)$. Then

$$\langle f, \phi_1 \rangle = 2, \quad \langle f, \phi_2 \rangle = -1, \quad \langle f, \phi_3 \rangle = 3.$$

Hence

$$f \sim 2\phi_1 - \phi_2 + 3\phi_3.$$

The following theorems show that the partial sums of the Fourier series of f have a certain minimum property. We shall assume here and in the rest of this chapter that $f \in \mathcal{R}$, although this hypothesis can be weakened.

Proposition 17.27. *Let (ϕ_n) be an orthonormal system of functions on $[a, b]$. Let*

$$s_n(x) = \sum_{k=1}^n c_k \phi_k(x)$$

be the n -th partial sum of the Fourier series of f , and let

$$t_n(x) = \sum_{k=1}^n \gamma_k \phi_k(x).$$

Then

$$\|f - s_n\| \leq \|f - t_n\|, \tag{17.8}$$

where equality holds if and only if $\gamma_k = c_k$ for $k = 1, \dots, n$.

That is to say, among all functions t_n , s_n gives the best possible mean square approximation to f .

Proof. We want to show that

$$\langle f - s_n, f - s_n \rangle \leq \langle f - t_n, f - t_n \rangle.$$

Note that

$$\begin{aligned} \langle f, s_n \rangle &= \left\langle f, \sum_{k=1}^n c_k \phi_k \right\rangle = \sum_{k=1}^n \overline{c_k} \langle f, \phi_k \rangle = \sum_{k=1}^n \overline{c_k} c_k = \sum_{k=1}^n |c_k|^2 \\ \langle s_n, s_n \rangle &= \left\langle \sum_{k=1}^n c_k \phi_k, \sum_{k=1}^n c_k \phi_k \right\rangle = \sum_{k=1}^n \langle c_k \phi_k, c_k \phi_k \rangle = \sum_{k=1}^n |c_k|^2 \\ \langle f, t_n \rangle &= \sum_{k=1}^n c_k \overline{\gamma_k} \\ \langle t_n, f \rangle &= \sum_{k=1}^n \gamma_k \overline{c_k} \\ \langle t_n, t_n \rangle &= \sum_{k=1}^n |\gamma_k|^2 \end{aligned}$$

Hence we rewrite the desired inequality as

$$\begin{aligned} &\iff \langle f, f \rangle - \sum_{k=1}^n |c_k|^2 \leq \langle f, f \rangle - \sum_{k=1}^n c_k \overline{\gamma_k} - \sum_{k=1}^n \gamma_k \overline{c_k} + \sum_{k=1}^n |\gamma_k|^2 \\ &\iff \sum_{k=1}^n (c_k \overline{c_k} - c_k \overline{\gamma_k} - \gamma_k \overline{c_k} + \gamma_k \overline{\gamma_k}) \geq 0 \\ &\iff \sum_{k=1}^n (c_k - \gamma_k)(\overline{c_k} - \overline{\gamma_k}) \geq 0 \\ &\iff \sum_{k=1}^n |c_k - \gamma_k|^2 \geq 0 \end{aligned}$$

which holds true. Then equality holds if and only if $|c_k - \gamma_k| = 0$, i.e.,

$$\gamma_k = c_k \quad (k = 1, \dots, n).$$

□

Proposition 17.28 (Bessel inequality). *Let (ϕ_n) be an orthonormal system of functions on $[a, b]$, and*

$$f(x) \sim \sum_{n=1}^{\infty} c_n \phi_n(x).$$

Then

$$\sum_{n=1}^{\infty} |c_n|^2 \leq \|f\|^2. \tag{17.9}$$

In particular, $c_n \rightarrow 0$.

Proof. Letting $n \rightarrow \infty$ in (72), we obtain (73)

□

the case where equality holds is called Parseval's identity

From now on we shall deal only with the trigonometric system. We shall consider functions f that have period 2π , and are Riemann-integrable on $[-\pi, \pi]$ (and hence on every bounded interval). The Fourier series of f is then the series (63) whose coefficients c_n are given by the integrals (62), and

$$s_N(x) = s_N(f; x) = \sum_{n=-N}^N c_n e^{inx}$$

is the N -th partial sum of the Fourier series of f . The inequality (72) now takes the form

In order to obtain an expression for s_N that is more manageable than (75) we introduce the *Dirichlet kernel*

$$D_N(x) := \sum_{n=-N}^N e^{inx}.$$

It follows that

$$\begin{aligned}
 D_N(x) &= \sum_{n=-N}^N e^{inx} \\
 &= \frac{e^{-iNx} \left[(e^{ix})^{2N+1} - 1 \right]}{e^{ix} - 1} \\
 &= \frac{e^{i(N+1)x} - e^{iNx}}{e^{ix} - 1} \\
 &= \frac{e^{i(N+\frac{1}{2})x} - e^{-i(N+\frac{1}{2})x}}{e^{\frac{ix}{2}} - e^{-\frac{ix}{2}}} \\
 &= \frac{\sin \left(N + \frac{1}{2} \right) x}{\sin \frac{1}{2} x}
 \end{aligned}$$

Then, for some dummy variable t ,

$$\begin{aligned}
 s_N(x) &= \sum_{n=-N}^N c_n e^{inx} = \sum_{n=-N}^N \left[\frac{1}{2\pi} \int_{-\pi}^{\pi} f(t) e^{-int} dt \right] e^{inx} \\
 &= \frac{1}{2\pi} \int_{-\pi}^{\pi} \left[\sum_{n=-N}^N f(t) e^{in(x-t)} \right] dt \\
 &= \frac{1}{2\pi} \int_{-\pi}^{\pi} f(t) \left[\sum_{n=-N}^N e^{in(x-t)} \right] dt \\
 &= \frac{1}{2\pi} \int_{-\pi}^{\pi} f(t) D_N(x-t) dt.
 \end{aligned}$$

Define the *convolution* of f and g as

$$(f * g)(t) := \int_E f(t)g(x-t) dt.$$

The periodicity of all functions involved shows that it is immaterial over which interval we integrate, as long as its length is 2π . This shows that

$$s_N(x) = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(t) D_N(x-t) dt = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x-t) D_N(t) dt.$$

We shall prove just one result about the pointwise convergence of Fourier series. Before that, we require the following result.

Proposition 17.29 (Riemann–Lebesgue lemma). *Let $f \in \mathcal{R}[a, b]$. Then*

$$\lim_{n \rightarrow \infty} \int_a^b f(x) \sin nx \, dx = 0. \quad (17.10)$$

Proof.

□

Proposition 17.30 (Pointwise convergence of Fourier series). *Suppose for some $x \in [-\pi, \pi]$ there exists $M > 0, \delta > 0$ such that*

$$\forall t \in (-\delta, \delta), \quad |f(x+t) - f(x)| \leq M|t|.$$

Then

$$\lim_{N \rightarrow \infty} s_N(f; x) = f(x).$$

Proof. Since

$$\frac{1}{2\pi} \int_{-\pi}^{\pi} D_N(x) dx = 1,$$

we can write

$$f(x) = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) D_N(t) dt.$$

Then

$$\begin{aligned} s_N(x) - f(x) &= \frac{1}{2\pi} \int_{-\pi}^{\pi} [f(x-t) - f(x)] D_N(t) dt \\ &= \frac{1}{2\pi} \int_{-\pi}^{\pi} [f(x-t) - f(x)] \frac{\sin\left(N + \frac{1}{2}\right)t}{\sin \frac{1}{2}t} dt \end{aligned}$$

Let $g(t) = \frac{f(x-t) - f(x)}{\sin \frac{1}{2}t}$, then

$$s_N(x) - f(x) = \frac{1}{2\pi} \int_{-\pi}^{\pi} g(t) \sin\left(N + \frac{1}{2}\right)t dt.$$

By the Riemann–Lebesgue lemma, we are done. \square

Corollary 17.31.

Here is another formulation of this corollary:

This is usually called the localisation theorem. It shows that the behaviour of the sequence $(s_N(f; x))$, as far as convergence is concerned, depends only on the values of f in some (arbitrarily small) neighbourhood of x . Two Fourier series may thus have the same behavior in one interval, but may behave in entirely different ways in some other interval. We have here a very striking contrast between Fourier series and power series (Theorem 8.5).

We conclude with two other approximation theorems.

Theorem 17.32. *If f is continuous (with period 2π) and if $\varepsilon > 0$, then there exists a trigonometric polynomial P such that*

$$|P(x) - f(x)| < \varepsilon \quad (x \in \mathbb{R}).$$

Proof. \square

Theorem 17.33 (Parseval's theorem). Suppose f and g are Riemann-integrable functions with period 2π , and

$$f(x) \sim \sum_{n=-\infty}^{\infty} c_n e^{inx}, \quad g(x) \sim \sum_{n=-\infty}^{\infty} \gamma_n e^{inx}.$$

Then

(i)

$$\lim_{N \rightarrow \infty} \|f - s_N(f)\|^2 = 0.$$

(ii)

$$\frac{1}{2\pi} \langle f, g \rangle = \sum_{n=-\infty}^{\infty} c_n \overline{\gamma_n}.$$

(iii)

$$\|f\|^2 = \sum_{n=-\infty}^{\infty} |c_n|^2.$$

Proof.

(i)

(ii)

(iii)

□

§17.4 Gamma Function

The *Gamma function* simulates the factorial.

Definition 17.34 (Gamma function). For $0 < x < \infty$, the *Gamma function* is defined as

$$\Gamma(x) := \int_0^{\infty} t^{x-1} e^{-t} dt. \quad (17.11)$$

The integral converges for these x . (When $x < 1$, both 0 and ∞ have to be looked at.)

Lemma 17.35.

(i) *The functional equation*

$$\Gamma(x+1) = x\Gamma(x)$$

holds for $0 < x < \infty$.

(ii) $\Gamma(n+1) = n!$ for $n = 1, 2, 3, \dots$

(iii) $\log \Gamma$ is convex on $(0, \infty)$.

Proof.

(i) Integrate by parts:

$$\begin{aligned} \Gamma(x+1) &= \int_0^{\infty} t^x e^{-t} dt \\ &= \left[-t^x e^{-t} \right]_0^{\infty} + \int_0^{\infty} x t^{x-1} e^{-t} dt \\ &= 0 + x\Gamma(x) = x\Gamma(x). \end{aligned}$$

(ii) We have

$$\Gamma(1) = \int_0^{\infty} e^{-t} dt = \left[-e^{-t} \right]_0^{\infty} = 1.$$

Since $\Gamma(1) = 1$, (i) implies (ii) by induction.

(iii) To show that $\log \Gamma(x)$ is convex, we need to show that for all $p, q > 1$ with $\frac{1}{p} + \frac{1}{q} = 1$,

$$\log \Gamma\left(\frac{x}{p} + \frac{y}{q}\right) \geq \frac{1}{p} \log \Gamma(x) + \frac{1}{q} \log \Gamma(y).$$

This is equivalent to showing

$$\Gamma\left(\frac{x}{p} + \frac{y}{q}\right) \geq \Gamma(x)^{\frac{1}{p}} + \Gamma(y)^{\frac{1}{q}}.$$

We have

$$\begin{aligned}
 \Gamma\left(\frac{x}{p} + \frac{y}{q}\right) &= \int_0^\infty t^{\frac{x}{p} + \frac{y}{q} - 1} e^{-t} dt \\
 &= \int_0^\infty t^{\frac{x-1}{p} + \frac{y-1}{q}} + e^{-t\left(\frac{1}{p} + \frac{1}{q}\right)} dt \\
 &= \int_0^\infty \left(t^{x-1} e^{-t}\right)^{\frac{1}{p}} \left(t^{y-1} e^{-t}\right)^{\frac{1}{q}} dt \\
 &\leq \left[\int_0^\infty \left(t^{\frac{x-1}{p}} e^{-\frac{t}{p}}\right)^p dt\right]^{\frac{1}{p}} \left[\int_0^\infty \left(t^{\frac{y-1}{q}} e^{-\frac{t}{q}}\right)^q dt\right]^{\frac{1}{q}} \\
 &= \Gamma(x)^{\frac{1}{p}} \Gamma(y)^{\frac{1}{q}}
 \end{aligned}$$

where the penultimate line holds as a result of Holder's inequality.

□

In fact, these three properties characterise Γ completely.

Lemma 17.36 (Characterisation of Γ). *If f is a positive function on $(0, \infty)$ such that*

(i) $f(x + 1) = xf(x)$,

(ii) $f(1) = 1$,

(iii) $\log f$ is convex,

then $f(x) = \Gamma(x)$.

Proof.

□

Definition 17.37 (Beta function). For $x > 0$ and $y > 0$, the *beta function* is defined as

$$B(x, y) := \int_0^1 t^{x-1} (1-t)^{y-1} dt.$$

Lemma 17.38.

$$B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}.$$

Proof. Let $f(x) = \frac{\Gamma(x+y)}{\Gamma(y)} B(x, y)$. We want to prove that $f(x) = \Gamma(x)$, using Lemma 17.36.

(i)

$$B(x+1, y) = \int_0^1 t^x (1-t)^{y-1} dt.$$

Integrating by parts gives

$$\begin{aligned} B(x+1, y) &= \underbrace{\left[t^x \cdot \frac{(1-t)^y}{y} (-1) \right]_0^1}_0 + \int_0^1 x t^{x-1} \frac{(1-t)^y}{y} dt \\ &= \frac{x}{y} \int_0^1 t^{x-1} (1-t)^{y-1} (1-t) dt \\ &= \frac{x}{y} \left(\int_0^1 t^{x-1} (1-t)^{y-1} dt - \int_0^1 t^x (1-t)^{y-1} dt \right) \\ &= \frac{x}{y} (B(x, y) - B(x+1, y)) \end{aligned}$$

which gives $B(x+1, y) = \frac{x}{x+y} B(x, y)$. Thus

$$\begin{aligned} f(x+1) &= \frac{\Gamma(x+1+y)}{\Gamma(y)} B(x+1, y) \\ &= \frac{(x+y)B(x+y)}{\Gamma(y)} \cdot \frac{x}{x+y} B(x, y) \\ &= x f(x). \end{aligned}$$

(ii)

$$B(1, y) = \int_0^1 (1-t)^{y-1} dt = \left[-\frac{(1-t)^y}{y} \right]_0^1 = \frac{1}{y}$$

and thus

$$f(1) = \frac{\Gamma(1+y)}{\Gamma(y)} B(1, y) = \frac{y\Gamma(y)}{\Gamma(y)} \frac{1}{y} = 1.$$

(iii) We now show that $\log B(x, y)$ is convex, so that

$$\log f(x) = \underbrace{\log \Gamma(x+y)}_{\text{convex}} + \log B(x, y) - \underbrace{\log \Gamma(y)}_{\text{constant}}$$

is convex with respect to x .

$$B(x_1, y)^{\frac{1}{p}} B(x_2, y)^{\frac{1}{q}} = \left(\int_0^1 t^{x_1-1} (1-t)^{y-1} dt \right)^{\frac{1}{p}} \left(\int_0^1 t^{x_2-1} (1-t)^{y-1} dt \right)^{\frac{1}{q}}$$

By Hölder's inequality,

$$\begin{aligned} B(x_1, y)^{\frac{1}{p}} B(x_2, y)^{\frac{1}{q}} &= \int_0^1 \left[t^{x_1-1} (1-t)^{y-1} \right]^{\frac{1}{p}} \left[t^{x_2-1} (1-t)^{y-1} \right]^{\frac{1}{q}} dt \\ &= \int_0^1 t^{\frac{x_1}{p} + \frac{x_2}{q} - 1} (1-t)^{y-1} dt \\ &= B\left(\frac{x_1}{p} + \frac{x_2}{q}, y\right). \end{aligned}$$

Taking log on both sides gives

$$\log B(x, y)^{\frac{1}{p}} B(x_2, y)^{\frac{1}{q}} \geq \log B\left(\frac{x_1}{p} + \frac{x_2}{q}, y\right)$$

or

$$\frac{1}{p} \log B(x, y) + \frac{1}{q} \log B(x_2, y) \geq \log B\left(\frac{x_1}{p} + \frac{x_2}{q}, y\right).$$

Hence $\log B(x, y)$ is convex, so $\log f(x, y)$ is convex.

Therefore $f(x, y) = \Gamma(x)\Gamma(y)$ which implies $B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}$. □

An alternative form of Γ is as follows:

$$\Gamma(x) = 2 \int_0^{+\infty} t^{2x-1} e^{-t^2} dt.$$

Using this form of Γ , we present an alternative proof.

Proof.

$$\begin{aligned} \Gamma(x)\Gamma(y) &= \left(2 \int_0^{+\infty} t^{2x-1} e^{-t^2} dt\right) \left(2 \int_0^{+\infty} s^{2y-1} e^{-s^2} ds\right) \\ &= 4 \iint_{[0,+\infty) \times [0,+\infty)} t^{2x-1} s^{2y-1} e^{-(t^2+s^2)} dt ds \end{aligned}$$

Using polar coordinates transformation, let $t = r \cos \theta$, $s = r \sin \theta$. Then $dt ds = r dr d\theta$. Thus

$$\begin{aligned} \Gamma(x)\Gamma(y) &= 4 \int_0^{\frac{\pi}{2}} \left[\int_0^{+\infty} r^{2x-1} \cos^{2x-1} \theta \cdot r^{2y-1} \sin^{2y-1} \theta \cdot e^{-r^2} \cdot r dr \right] d\theta \\ &= 2 \underbrace{\int_0^{\frac{\pi}{2}} \cos^{2x-1} \theta \sin^{2y-1} \theta d\theta}_{B(x,y)} \cdot 2 \underbrace{\int_0^{+\infty} r^{2(x+y)-1} e^{-r^2} dr}_{\Gamma(x+y)} \end{aligned}$$

since

$$\begin{aligned} B(x, y) &= \int_0^1 t^{x-1} (1-t)^{y-1} dt \quad t = \cos^2 \theta \\ &= \int_{\frac{\pi}{2}}^0 \cos^{2(x-1)} \theta \sin^{2(y-1)} \theta \cdot 2 \cos \theta (-\sin \theta) d\theta \\ &= 2 \int_0^{\frac{\pi}{2}} \cos^{2x-1} \theta \sin^{2y-1} \theta d\theta. \end{aligned}$$

Hence $B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}$. □

More on polar coordinates:

$$I = \int_{-\infty}^{+\infty} e^{-x^2} dx = \sqrt{\pi}. \tag{17.12}$$

Proof.

$$\begin{aligned}
 I^2 &= \int_{-\infty}^{+\infty} e^{-x^2} dx \int_{-\infty}^{+\infty} e^{-y^2} dy \\
 &= \iint_{\mathbb{R}^2} e^{-(x^2+y^2)} dx dy \quad x = r \cos \theta, y = r \sin \theta \\
 &= \int_0^{2\pi} \underbrace{\int_0^{+\infty} e^{-r^2} r dr}_{\text{constant w.r.t. } \theta} d\theta \quad s = r^2, ds = 2r dr \\
 &= 2\pi \int_0^{+\infty} e^{-s} \cdot \frac{1}{2} ds \\
 &= 2\pi \left[\frac{1}{2} e^{-s} (-1) \right]_0^{\infty} = \pi
 \end{aligned}$$

and thus

$$I = \int_{-\infty}^{+\infty} e^{-x^2} dx = \sqrt{\pi}.$$

□

From this, we have

$$\Gamma\left(\frac{1}{2}\right) = 2 \int_0^{\infty} e^{-t^2} dt = \sqrt{\pi}.$$

Lemma 17.39.

$$\Gamma(x) = \frac{2^{x-1}}{\sqrt{\pi}} \Gamma\left(\frac{x}{2}\right) \Gamma\left(\frac{x+1}{2}\right).$$

Proof. Let $f(x) = \frac{2^{x-1}}{\sqrt{\pi}} \Gamma\left(\frac{x}{2}\right) \Gamma\left(\frac{x+1}{2}\right)$. We want to prove that $f(x) = \Gamma(x)$.

(i)

$$\begin{aligned}
 f(x+1) &= \frac{2^x}{\sqrt{\pi}} \Gamma\left(\frac{x+1}{2}\right) \Gamma\left(\frac{x}{2} + 1\right) \\
 &= \frac{2^x}{\sqrt{\pi}} \Gamma\left(\frac{x+1}{2}\right) \frac{x}{2} \Gamma\left(\frac{x}{2}\right) \\
 &= x f(x)
 \end{aligned}$$

(ii) $f(1) = \frac{1}{\sqrt{\pi}} \Gamma\left(\frac{1}{2}\right) \Gamma(1) = 1$ since $\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$.

(iii)

$$\log f(x) = \underbrace{(x-1) \log 2}_{\text{linear}} + \underbrace{\log \Gamma\left(\frac{x}{2}\right)}_{\text{convex}} + \underbrace{\log \Gamma\left(\frac{x+1}{2}\right)}_{\text{convex}} - \underbrace{\log \sqrt{\pi}}_{\text{constant}}$$

and hence $\log f(x)$ is convex.

Therefore $f(x) = \Gamma(x)$.

□

Theorem 17.40 (Stirling's formula). *This provides a simple approximate expression for $\Gamma(x + 1)$ when x is large (hence for $n!$ when n is large). The formula is*

$$\lim_{x \rightarrow \infty} \frac{\Gamma(x + 1)}{(x/e)^x \sqrt{2\pi x}} = 1. \quad (17.13)$$

Proof. □

Lemma 17.41.

$$B(p, 1 - p) = \Gamma(p)\Gamma(1 - p) = \frac{\pi}{\sin p\pi}.$$

Proof. We have

$$\begin{aligned} B(p, 1 - p) &= \int_0^1 t^{p-1}(1-t)^{-p} dt \\ &= \int_0^\infty \left(\frac{x}{1+x}\right)^{p-1} \left(\frac{1}{1+x}\right)^{-p} \frac{1}{(1+x)^2} dx \quad [x = \frac{t}{1-t}] \\ &= \int_0^\infty \frac{x^{p-1}}{1+x} dx \\ &= \int_0^1 \frac{x^{p-1}}{1+x} dx + \int_1^\infty \frac{x^{p-1}}{1+x} dx \end{aligned}$$

See that

$$\begin{aligned} \int_1^\infty \frac{x^{p-1}}{1+x} dx &= \int_1^0 \frac{y^{1-p}}{1+\frac{1}{y}} \left(-\frac{1}{y^2}\right) dy \quad [x = \frac{1}{y}] \\ &= \int_0^1 \frac{y^{-p}}{1+y} dy = \int_0^1 \frac{x^{-p}}{1+x} dx \end{aligned}$$

so

$$\begin{aligned}
B(p, 1-p) &= \int_0^1 \frac{x^{p-1} + x^{-p}}{1+x} dx \\
&= \lim_{r \rightarrow 1^-} \int_0^r (x^{p-1} + x^{-p}) \sum_{k=0}^{\infty} (-1)^k x^k dx \\
&= \lim_{r \rightarrow 1^-} \int_0^r \left(\sum_{k=0}^{\infty} (-1)^k x^{k+p-1} + \sum_{k=0}^{\infty} (-1)^k x^{k-p} \right) dx \\
&= \lim_{r \rightarrow 1^-} \left[\sum_{k=0}^{\infty} (-1)^k \frac{x^{k+p}}{k+p} + \sum_{k=0}^{\infty} (-1)^k \frac{x^{k-p+1}}{k-p+1} \right]_0^r \\
&= \sum_{k=0}^{\infty} (-1)^k \frac{1}{k+p} + \sum_{k=0}^{\infty} (-1)^k \frac{1}{k-p+1} \\
&= \frac{1}{p} + \sum_{k=1}^{\infty} (-1)^k \frac{1}{k+p} + \sum_{k=1}^{\infty} (-1)^{k-1} \frac{1}{k+p} \\
&= \frac{1}{p} + \sum_{k=1}^{\infty} \frac{(-1)^k 2p}{p^2 - k^2}
\end{aligned}$$

□

Exercises

V

Multivariable Analysis

Vector functions, limits and continuity, vector derivatives, total derivative, chain rule, partial derivatives, gradient, inverse function theorem, implicit function theorem.

Multiple integrals, Fubini theorem, change of variables.

Differential forms, closed and exact forms, wedge product, simplexes and chains, integration on chains, manifolds, integration on manifolds, Stokes Theorem. The volume element, connections with the classical theorems of vector calculus.

18 Functions of Several Variables

We shall now switch to a different topic, namely that of differentiation in several variable calculus. More precisely, we shall be dealing with maps $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ from one Euclidean space to another, and trying to understand what the derivative of such a map is.

Before we do so, however, we need to recall some notions from linear algebra, most importantly that of a linear map and a matrix.

§18.1 Linear Transformations

Some linear algebra prerequisites.

§18.2 Differentiation

The Derivative

Recall that for $f: \mathbb{R} \rightarrow \mathbb{R}$, we defined the derivative at x as

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}.$$

In other words, there exists a number a (the derivative of f at x) such that

$$\lim_{h \rightarrow 0} \left| \frac{f(x+h) - f(x)}{h} - a \right| = \lim_{h \rightarrow 0} \left| \frac{f(x+h) - f(x) - ah}{h} \right| = \lim_{h \rightarrow 0} \frac{|f(x+h) - f(x) - ah|}{|h|} = 0.$$

Multiplying by a is a linear map in one dimension: $h \mapsto ah$. Namely, we think of $a \in \mathcal{L}(\mathbb{R}, \mathbb{R})$. Hence we can use this definition to extend differentiation to more variables.

Definition 18.1 (Derivative). Let $U \subset \mathbb{R}^n$ be open, $f: U \rightarrow \mathbb{R}^m$. We say f is *differentiable* at $\mathbf{x} \in U$ if there exists $A \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m)$ such that

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}} \frac{\|f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x}) - A\mathbf{h}\|}{\|\mathbf{h}\|} = 0.$$

Then we write $f'(\mathbf{x}) = A$, and say that A is the *derivative* of f at \mathbf{x} .

If f is differentiable at every $\mathbf{x} \in U$, we say f is *differentiable on U* .

Remark. Note that the derivative is a function from U to $\mathcal{L}(\mathbb{R}^n, \mathbb{R}^m)$.

We now show that the above derivative is in fact unique.

Lemma 18.2 (Uniqueness of derivative). *Let $U \subset \mathbb{R}^n$ be open, $\mathbf{f}: U \rightarrow \mathbb{R}^m$. Suppose $\mathbf{x} \in U$, and there exist $A, B \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m)$ such that*

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}} \frac{\|\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) - A\mathbf{h}\|}{\|\mathbf{h}\|} = 0 \quad \text{and} \quad \lim_{\mathbf{h} \rightarrow \mathbf{0}} \frac{\|\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) - B\mathbf{h}\|}{\|\mathbf{h}\|} = 0.$$

Then $A = B$.

Proof. Suppose $\mathbf{h} \neq \mathbf{0}$. Compute

$$\begin{aligned} \frac{\|(A - B)\mathbf{h}\|}{\|\mathbf{h}\|} &= \frac{\|(\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) - A\mathbf{h}) - (\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) - B\mathbf{h})\|}{\|\mathbf{h}\|} \\ &\leq \frac{\|\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) - A\mathbf{h}\|}{\|\mathbf{h}\|} + \frac{\|\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) - B\mathbf{h}\|}{\|\mathbf{h}\|}. \end{aligned}$$

Taking the limit $\mathbf{h} \rightarrow \mathbf{0}$ on both sides gives

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}} \frac{\|(A - B)\mathbf{h}\|}{\|\mathbf{h}\|} = 0.$$

Thus given $\varepsilon > 0$, for all non-zero \mathbf{h} in some δ -ball around the origin we have

$$\frac{\|(A - B)\mathbf{h}\|}{\|\mathbf{h}\|} = \left\| (A - B) \frac{\mathbf{h}}{\|\mathbf{h}\|} \right\| < \varepsilon.$$

For any given $\mathbf{v} \in \mathbb{R}^n$ with $\|\mathbf{v}\| = 1$, if $\mathbf{h} = \frac{\delta}{2}\mathbf{v}$, then $\|\mathbf{h}\| < \delta$ and $\frac{\mathbf{h}}{\|\mathbf{h}\|} = \mathbf{v}$. So $\|(A - B)\mathbf{v}\| < \varepsilon$. Taking the supremum over all \mathbf{v} with $\|\mathbf{v}\| = 1$, we get the operator norm $\|A - B\| \leq \varepsilon$.

As $\varepsilon > 0$ was arbitrary, we must have $\|A - B\| = 0$, or in other words $A = B$. \square

Example 18.3. If $\mathbf{f}(\mathbf{x}) = A\mathbf{x}$ for a linear mapping A , then $\mathbf{f}'(\mathbf{x}) = A$:

$$\frac{\|\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) - A\mathbf{h}\|}{\|\mathbf{h}\|} = \frac{\|A(\mathbf{x} + \mathbf{h}) - A\mathbf{x} - A\mathbf{h}\|}{\|\mathbf{h}\|} = \frac{0}{\|\mathbf{h}\|} = 0.$$

Lemma 18.4 (Differentiability implies continuity). *Let $U \subset \mathbb{R}^n$ be open, suppose $\mathbf{f}: U \rightarrow \mathbb{R}^m$ is differentiable at $\mathbf{x} \in U$. Then \mathbf{f} is continuous at \mathbf{x} .*

Proof. Another way to write the differentiability of \mathbf{f} at \mathbf{x} is to consider the *remainder*:

$$\mathbf{r}(\mathbf{h}) := \mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) - \mathbf{f}'(\mathbf{x})\mathbf{h}.$$

By definition, \mathbf{f} is differentiable at \mathbf{x} if $\lim_{\mathbf{h} \rightarrow \mathbf{0}} \frac{\|\mathbf{r}(\mathbf{h})\|}{\|\mathbf{h}\|} = 0$, so $\mathbf{r}(\mathbf{h})$ itself goes to zero.

The mapping $\mathbf{h} \mapsto \mathbf{f}'(\mathbf{x})\mathbf{h}$ is a linear mapping between finite-dimensional spaces, hence continuous and $\lim_{\mathbf{h} \rightarrow \mathbf{0}} \mathbf{f}'(\mathbf{x})\mathbf{h} = \mathbf{0}$. Thus, $\lim_{\mathbf{h} \rightarrow \mathbf{0}} \mathbf{f}(\mathbf{x} + \mathbf{h}) = \mathbf{f}(\mathbf{x})$. That is, \mathbf{f} is continuous at \mathbf{x} . \square

Differentiation is a linear map on the space of differentiable functions.

Lemma 18.5. Let $U \subset \mathbb{R}^n$ be open, suppose $\mathbf{f}, \mathbf{g} : U \rightarrow \mathbb{R}^m$ are differentiable at $\mathbf{x} \in U$, let $\alpha \in \mathbb{R}$.
Then

(i) $\mathbf{f} + \mathbf{g}$ is differentiable at \mathbf{x} , and (addition)

$$(\mathbf{f} + \mathbf{g})'(\mathbf{x}) = \mathbf{f}'(\mathbf{x}) + \mathbf{g}'(\mathbf{x}).$$

(ii) $\alpha\mathbf{f}$ is differentiable at \mathbf{x} , and (scalar multiplication)

$$(\alpha\mathbf{f})'(\mathbf{x}) = \alpha\mathbf{f}'(\mathbf{x}).$$

Proof. Let $\mathbf{h} \in \mathbb{R}^n$, $h \neq \mathbf{0}$.

(i) Write

$$\begin{aligned} & \frac{\|\mathbf{f}(\mathbf{x} + \mathbf{h}) + \mathbf{g}(\mathbf{x} + \mathbf{h}) - (\mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})) - (\mathbf{f}'(\mathbf{x}) + \mathbf{g}'(\mathbf{x}))\mathbf{h}\|}{\|\mathbf{h}\|} \\ & \leq \frac{\|\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) - \mathbf{f}'(\mathbf{x})\mathbf{h}\|}{\|\mathbf{h}\|} + \frac{\|\mathbf{g}(\mathbf{x} + \mathbf{h}) - \mathbf{g}(\mathbf{x}) - \mathbf{g}'(\mathbf{x})\mathbf{h}\|}{\|\mathbf{h}\|} \end{aligned}$$

Then take limits $\mathbf{h} \rightarrow \mathbf{0}$ on both sides of the equation.

(ii) Write

$$\frac{\|\alpha\mathbf{f}(\mathbf{x} + \mathbf{h}) - \alpha\mathbf{f}(\mathbf{x}) - \alpha\mathbf{f}'(\mathbf{x})\mathbf{h}\|}{\|\mathbf{h}\|} = |\alpha| \frac{\|\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) - \mathbf{f}'(\mathbf{x})\mathbf{h}\|}{\|\mathbf{h}\|}.$$

Then take limits $\mathbf{h} \rightarrow \mathbf{0}$ on both sides of the equation.

□

We now extend the chain rule to the present situation.

Lemma 18.6 (Chain rule). Let $U \subset \mathbb{R}^n$, $V \subset \mathbb{R}^m$ be open. Suppose $\mathbf{f} : U \rightarrow \mathbb{R}^m$ is differentiable at $\mathbf{x} \in U$, $\mathbf{f}(U) \subset V$, and $g : V \rightarrow \mathbb{R}^k$ is differentiable at $\mathbf{f}(\mathbf{x})$.
Then $\mathbf{F} = g \circ \mathbf{f}$ is differentiable at \mathbf{x} , and

$$\mathbf{F}'(\mathbf{x}) = g'(\mathbf{f}(\mathbf{x})) \mathbf{f}'(\mathbf{x}).$$

Without the points where things are evaluated, we write $\mathbf{F}' = (g \circ \mathbf{f})' = g'\mathbf{f}'$. The derivative of the composition $g \circ \mathbf{f}$ is the composition of the derivatives of g and \mathbf{f} : If $\mathbf{f}'(\mathbf{x}) = A$ and $g'(\mathbf{f}(\mathbf{x})) = B$, then $\mathbf{F}'(\mathbf{x}) = BA$, just as for linear maps.

Proof. Let $A = \mathbf{f}'(\mathbf{x})$ and $B = g'(\mathbf{f}(\mathbf{x}))$. Take a non-zero $\mathbf{h} \in \mathbb{R}^n$ and write $\mathbf{y} = \mathbf{f}(\mathbf{x})$, $\mathbf{k} = \mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x})$.
Let

$$\mathbf{r}(\mathbf{h}) = \mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) - A\mathbf{h}.$$

Then $\mathbf{r}(\mathbf{h}) = \mathbf{k} - A\mathbf{h}$ or $A\mathbf{h} = \mathbf{k} - \mathbf{r}(\mathbf{h})$, and $\mathbf{f}(\mathbf{x} + \mathbf{h}) = \mathbf{y} + \mathbf{k}$. We look at the quantity we need to go to zero:

$$\begin{aligned} \frac{\|\mathbf{F}(\mathbf{x} + \mathbf{h}) - \mathbf{F}(\mathbf{x}) - BA\mathbf{h}\|}{\|\mathbf{h}\|} &= \frac{\|\mathbf{g}(\mathbf{f}(\mathbf{x} + \mathbf{h})) - \mathbf{g}(\mathbf{f}(\mathbf{x})) - BA\mathbf{h}\|}{\|\mathbf{h}\|} \\ &= \frac{\|\mathbf{g}(\mathbf{y} + \mathbf{k}) - \mathbf{g}(\mathbf{y}) - B(\mathbf{k} - \mathbf{r}(\mathbf{h}))\|}{\|\mathbf{h}\|} \\ &\leq \frac{\|\mathbf{g}(\mathbf{y} + \mathbf{k}) - \mathbf{g}(\mathbf{y}) - B\mathbf{k}\|}{\|\mathbf{h}\|} + \|B\| \frac{\|\mathbf{r}(\mathbf{h})\|}{\|\mathbf{h}\|} \\ &= \frac{\|\mathbf{g}(\mathbf{y} + \mathbf{k}) - \mathbf{g}(\mathbf{y}) - B\mathbf{k}\|}{\|\mathbf{k}\|} \frac{\|\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x})\|}{\|\mathbf{h}\|} + \|B\| \frac{\|\mathbf{r}(\mathbf{h})\|}{\|\mathbf{h}\|}. \end{aligned}$$

Take the limit $\mathbf{h} \rightarrow \mathbf{0}$. We examine the three terms:

- Since \mathbf{f} is differentiable at \mathbf{x} , $\lim_{\mathbf{h} \rightarrow \mathbf{0}} \frac{\|\mathbf{r}(\mathbf{h})\|}{\|\mathbf{h}\|} = 0$.
- Since \mathbf{f} is continuous at \mathbf{x} , $\mathbf{k} \rightarrow \mathbf{0}$ as $\mathbf{h} \rightarrow \mathbf{0}$. Thus since \mathbf{g} is differentiable at \mathbf{y} ,

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}} \frac{\|\mathbf{g}(\mathbf{y} + \mathbf{k}) - \mathbf{g}(\mathbf{y}) - B\mathbf{k}\|}{\|\mathbf{k}\|} = 0.$$

- Write

$$\begin{aligned} \frac{\|\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x})\|}{\|\mathbf{h}\|} &\leq \frac{\|\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) - A\mathbf{h}\|}{\|\mathbf{h}\|} + \frac{\|A\mathbf{h}\|}{\|\mathbf{h}\|} \\ &\leq \frac{\|\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) - A\mathbf{h}\|}{\|\mathbf{h}\|} + \|A\|. \end{aligned}$$

Since \mathbf{f} is differentiable at \mathbf{x} , for small enough \mathbf{h} , the quantity $\frac{\|\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) - A\mathbf{h}\|}{\|\mathbf{h}\|}$ is bounded. Thus the term $\frac{\|\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x})\|}{\|\mathbf{h}\|}$ stays bounded as $\mathbf{h} \rightarrow \mathbf{0}$.

Therefore

$$\lim_{\mathbf{h} \rightarrow \mathbf{0}} \frac{\|\mathbf{F}(\mathbf{x} + \mathbf{h}) - \mathbf{F}(\mathbf{x}) - BA\mathbf{h}\|}{\|\mathbf{h}\|} = 0,$$

so $\mathbf{F}'(\mathbf{x}) = BA$ as desired. □

Partial Derivatives

There is another way to generalise the derivative from one dimension. We can simply hold all but one variables constant and take the regular derivative.

Let $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$ and $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$ be the standard bases of \mathbb{R}^n and \mathbb{R}^m .

Definition 18.7 (Partial derivative). Let $U \subset \mathbb{R}^n$ be open, $\mathbf{f}: U \rightarrow \mathbb{R}^m$. The *components* of \mathbf{f} are the real functions f_1, \dots, f_m defined by

$$\mathbf{f}(\mathbf{x}) = \sum_{i=1}^m f_i(\mathbf{x})\mathbf{u}_i \quad (\mathbf{x} \in U).$$

The *partial derivative* at $\mathbf{x} \in U$ is defined as

$$\frac{\partial f_i}{\partial x_j}(\mathbf{x}) := \lim_{t \rightarrow 0} \frac{f_i(\mathbf{x} + t\mathbf{e}_j) - f_i(\mathbf{x})}{t} \quad (1 \leq i \leq m, 1 \leq j \leq n),$$

provided the limit exists.

Partial derivatives are easier to compute with all the machinery of calculus, and they provide a way to compute the total derivative of a function.

Proposition 18.8. Let $U \subset \mathbb{R}^n$ be open, suppose $\mathbf{f}: U \rightarrow \mathbb{R}^m$ is differentiable at $\mathbf{x} \in U$. Then all the partial derivatives at \mathbf{x} exist; with respect to the standard bases of \mathbb{R}^n and \mathbb{R}^m ,

$$[\mathbf{f}'(\mathbf{x})] = \begin{pmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{x}) & \frac{\partial f_1}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial f_1}{\partial x_n}(\mathbf{x}) \\ \frac{\partial f_2}{\partial x_1}(\mathbf{x}) & \frac{\partial f_2}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial f_2}{\partial x_n}(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1}(\mathbf{x}) & \frac{\partial f_m}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial f_m}{\partial x_n}(\mathbf{x}) \end{pmatrix}.$$

In other words,

$$\mathbf{f}'(\mathbf{x})\mathbf{e}_j = \sum_{i=1}^m \frac{\partial f_i}{\partial x_j}(\mathbf{x})\mathbf{u}_i \quad (j = 1, \dots, n).$$

The matrix of $\mathbf{f}'(\mathbf{x})$ is often called the *Jacobian matrix*.

Proof. Fix j . Since \mathbf{f} is differentiable at \mathbf{x} ,

$$\mathbf{f}(\mathbf{x} + t\mathbf{e}_j) - \mathbf{f}(\mathbf{x}) = \mathbf{f}'(\mathbf{x})(t\mathbf{e}_j) + \mathbf{r}(t\mathbf{e}_j)$$

where $\frac{\|\mathbf{r}(t\mathbf{e}_j)\|}{t} \rightarrow 0$ as $t \rightarrow 0$. Taking the limit $t \rightarrow 0$ on both sides, the linearity of $\mathbf{f}'(\mathbf{x})$ shows that

$$\lim_{t \rightarrow 0} \frac{\mathbf{f}(\mathbf{x} + t\mathbf{e}_j) - \mathbf{f}(\mathbf{x})}{t} = \mathbf{f}'(\mathbf{x})\mathbf{e}_j.$$

If we now represent \mathbf{f} in terms of its components, the above equation becomes

$$\lim_{t \rightarrow 0} \sum_{i=1}^m \frac{f_i(\mathbf{x} + t\mathbf{e}_j) - f_i(\mathbf{x})}{t} \mathbf{u}_i = \mathbf{f}'(\mathbf{x})\mathbf{e}_j.$$

It follows that each quotient in this sum has a limit as $t \rightarrow 0$ (see Theorem 4.10), so that each partial derivative $\frac{\partial f_i}{\partial x_j}$ exists. □

Gradients, Curves, and Directional Derivatives

Let γ be a differentiable mapping of $(a, b) \subset \mathbb{R}$ into an open set $U \subset \mathbb{R}^n$; that is, γ is a differentiable curve in U . Let $f: U \rightarrow \mathbb{R}$ be differentiable.

For $t \in (a, b)$, define

$$g(t) = f(\gamma(t)).$$

By the chain rule,

$$g'(t) = f'(\gamma(t)) \gamma'(t).$$

Since $\gamma'(t) \in \mathcal{L}(\mathbb{R}, \mathbb{R}^n)$ and $f'(\gamma(t)) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R})$, $g'(t)$ is a linear operator on \mathbb{R} ; thus, we can regard $g'(t)$ as a real number. This number can be computed in terms of the partial derivatives of f and the derivatives of the components of γ , as we shall now see.

With respect to the standard basis $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$ of \mathbb{R}^n , the matrix of $\gamma'(t)$ is the $n \times 1$ matrix which has $\gamma'_i(t)$ in the i -th row, where $\gamma_1, \dots, \gamma_n$ are the components of γ . For every $\mathbf{x} \in U$, the matrix of $f'(\mathbf{x})$ is the $1 \times n$ matrix which has $\frac{\partial f}{\partial x_j}$ in the j -th column. Hence the matrix of $g'(t)$ is the 1×1 matrix whose only entry is the real number

$$g'(t) = f'(\gamma(t))\gamma'(t) = \sum_{j=1}^n \frac{\partial f}{\partial x_j}(\gamma(t)) \frac{d\gamma_j}{dt} = \sum_{j=1}^n \frac{\partial f}{\partial x_j} \frac{d\gamma_j}{dt}.$$

Definition 18.9 (Gradient). Let $U \subset \mathbb{R}^n$ be open, suppose $f: U \rightarrow \mathbb{R}$ is differentiable. The **gradient** at $\mathbf{x} \in U$ is defined as

$$(\nabla f)(\mathbf{x}) := \sum_{j=1}^n \frac{\partial f}{\partial x_j}(\mathbf{x}) \mathbf{e}_j. \quad (18.1)$$

Writing $\gamma'(t)$ as components

$$\gamma'(t) = \sum_{j=1}^n \gamma'_j(t) \mathbf{e}_j,$$

using the scalar product, we can rewrite $g'(t)$ as

$$g'(t) = (\nabla f)(\gamma(t)) \cdot \gamma'(t). \quad (18.2)$$

Let us now fix $\mathbf{x} \in U$, take a unit vector $\mathbf{u} \in \mathbb{R}^n$, and let γ be

$$\gamma(t) = \mathbf{x} + t\mathbf{u}.$$

Then $\gamma'(t) = \mathbf{u}$ for every t . Hence Eq. (18.2) shows that

$$g'(0) = (\nabla f)(\mathbf{x}) \cdot \mathbf{u}.$$

On the other hand, we have

$$g(t) - g(0) = f(\mathbf{x} + t\mathbf{u}) - f(\mathbf{x}).$$

Hence

$$\lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{u}) - f(\mathbf{x})}{t} = (\nabla f)(\mathbf{x}) \cdot \mathbf{u}. \quad (18.3)$$

We call this limit the **directional derivative** of f at \mathbf{x} , in the direction of the unit vector \mathbf{u} , and may be denoted

by $(D_{\mathbf{u}}f)(\mathbf{x})$.

If f and \mathbf{x} are fixed, but \mathbf{u} varies, then Eq. (18.3) shows that $(D_{\mathbf{u}}f)(\mathbf{x})$ attains its maximum when \mathbf{u} is a positive scalar multiple of $(\nabla f)(\mathbf{x})$. [The case $(\nabla f)(\mathbf{x}) = \mathbf{0}$ should be excluded here.]

If $\mathbf{u} = \sum_j u_j \mathbf{e}_j$, then Eq. (18.3) shows that $(D_{\mathbf{u}}f)(\mathbf{x})$ can be expressed in terms of the partial derivatives of f at \mathbf{x} :

$$(D_{\mathbf{u}}f)(\mathbf{x}) = \sum_{j=1}^n \frac{\partial f}{\partial x_j}(\mathbf{x}) \mathbf{u}_j. \quad (18.4)$$

Proposition 18.10. *Let $U \subset \mathbb{R}^n$ be open and convex, suppose $\mathbf{f}: U \rightarrow \mathbb{R}^m$ is differentiable on U , and there exists a real number M such that*

$$\|\mathbf{f}'(\mathbf{x})\| \leq M \quad (\mathbf{x} \in U).$$

Then for all $\mathbf{a}, \mathbf{b} \in U$,

$$|\mathbf{f}(\mathbf{b}) - \mathbf{f}(\mathbf{a})| \leq M|\mathbf{b} - \mathbf{a}|.$$

Proof. Fix $\mathbf{a}, \mathbf{b} \in U$. Define

$$\gamma(t) = (1-t)\mathbf{a} + t\mathbf{b}$$

for all $t \in \mathbb{R}$ such that $\gamma(t) \in U$. Since U is convex, $\gamma(t) \in U$ if $0 \leq t \leq 1$. Put

$$\mathbf{g}(t) = \mathbf{f}(\gamma(t)).$$

Then

$$\mathbf{g}'(t) = \mathbf{f}'(\gamma(t)) \gamma'(t) = \mathbf{f}'(\gamma(t)) (\mathbf{b} - \mathbf{a}),$$

so that

$$|\mathbf{g}'(t)| \leq \|\mathbf{f}'(\gamma(t))\| |\mathbf{b} - \mathbf{a}| \leq M|\mathbf{b} - \mathbf{a}|$$

for all $t \in [0, 1]$. By Theorem 5.19,

$$|\mathbf{g}(1) - \mathbf{g}(0)| \leq M|\mathbf{b} - \mathbf{a}|.$$

But $\mathbf{g}(0) = \mathbf{f}(\mathbf{a})$ and $\mathbf{g}(1) = \mathbf{f}(\mathbf{b})$. This completes the proof. \square

Corollary 18.11. *If, in addition, $\mathbf{f}'(\mathbf{x}) = \mathbf{0}$ for all $\mathbf{x} \in U$, then \mathbf{f} is constant.*

Proof. To prove this, note that the hypotheses of the previous result hold now with $M = 0$. \square

The Jacobian

Definition 18.12 (Jacobian). Let $U \subset \mathbb{R}^n$, suppose $\mathbf{f}: U \rightarrow \mathbb{R}^m$ is differentiable. Define the *Jacobian* of \mathbf{f} at $\mathbf{x} \in U$ as

$$J_{\mathbf{f}}(\mathbf{x}) := \det[\mathbf{f}'(\mathbf{x})].$$

We shall also denote $J_{\mathbf{f}}$ as

$$\frac{\partial(f_1, \dots, f_n)}{\partial(x_1, \dots, x_n)}.$$

This last piece of notation may seem somewhat confusing, but it is quite useful when we need to specify the exact variables and function components used, as we will do, for example, in the implicit function theorem.

The Jacobian determinant $J_{\mathbf{f}}$ is a real-valued function, and when $n = 1$ it is simply the derivative. From the chain rule and $\det AB = \det A \det B$, it follows that

$$J_{\mathbf{f} \circ \mathbf{g}}(\mathbf{x}) = J_{\mathbf{f}}(\mathbf{g}(\mathbf{x}))J_{\mathbf{g}}(\mathbf{x}).$$

The determinant of a linear mapping tells us what happens to area/volume under the mapping. Similarly, the Jacobian determinant measures how much a differentiable mapping stretches things locally, and if it flips orientation. In particular, if the Jacobian determinant is non-zero then we would assume that locally the mapping is invertible (and we would be correct as we will later see).

§18.3 Continuity and The Derivative

Let us prove a “mean value theorem” for vector-valued functions.

Theorem 18.13. *If $\phi : [a, b] \rightarrow \mathbb{R}^n$ is differentiable on (a, b) and continuous on $[a, b]$, then there exists $t \in [a, b]$ such that*

$$\|\phi(b) - \phi(a)\| \leq (b - a)\|\phi'(t)\|. \quad (18.5)$$

We say $\mathbf{f} : U \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ is **continuously differentiable** if \mathbf{f} is differentiable, and $\mathbf{f}' : U \rightarrow \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m)$ is continuous; we also say that \mathbf{f} is a \mathcal{C}' -mapping, or that $\mathbf{f} \in \mathcal{C}'(U)$.

Proposition 18.14. *Let $U \subset \mathbb{R}^n$ be open, $\mathbf{f} : U \rightarrow \mathbb{R}^m$. Then \mathbf{f} is continuously differentiable if and only if all the partial derivatives $\frac{\partial f_i}{\partial x_j}$ exist and are continuous on U .*

Proof.

\implies Suppose $\mathbf{f} \in \mathcal{C}'(U)$.

\impliedby

□

§18.4 Inverse and Implicit Function Theorems

Inverse Function Theorem

The inverse function theorem states, roughly speaking, that a continuously differentiable mapping \mathbf{f} is invertible in a neighbourhood of any point \mathbf{x} at which the linear map $\mathbf{f}'(\mathbf{x})$ is invertible:

Theorem 18.15 (Inverse function theorem). *Let $U \subset \mathbb{R}^n$ be open, suppose $\mathbf{f}: U \rightarrow \mathbb{R}^n$ is a C^1 -mapping, and $\mathbf{f}'(\mathbf{a})$ is invertible for some $\mathbf{a} \in U$, and $\mathbf{b} = \mathbf{f}(\mathbf{a})$. Then*

(i) *there exist open sets $U, V \subset \mathbb{R}^n$ such that $\mathbf{a} \in U$, $\mathbf{b} \in V$, \mathbf{f} is bijective on U , and $\mathbf{f}(U) = V$;*

(ii) *if \mathbf{g} is the inverse of \mathbf{f} [which exists, by (i)], defined in V by*

$$\mathbf{g}(\mathbf{f}(\mathbf{x})) = \mathbf{x} \quad (\mathbf{x} \in U),$$

then $\mathbf{g} \in C^1(V)$.

Corollary 18.16. *Let $U \subset \mathbb{R}^n$ be open,*

Implicit Function Theorem

§18.5 Derivatives of Higher Order

Definition 18.17. Let $U \subset \mathbb{R}^n$ be open, suppose $f: U \rightarrow \mathbb{R}$, with partial derivatives $\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n}$. If the functions $\frac{\partial f}{\partial x_j}$ are themselves differentiable, then the *second-order partial derivatives* of f are defined by

$$\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial}{\partial x_i} \left(\frac{\partial f}{\partial x_j} \right) \quad (i, j = 1, \dots, n).$$

If all these functions $\frac{\partial^2 f}{\partial x_i \partial x_j}$ are continuous on U , we say that f is of class \mathcal{C}'' in U , or that $f \in \mathcal{C}''(U)$. $\mathbf{f}: U \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ is said to be of class \mathcal{C}'' if each component of \mathbf{f} is of class \mathcal{C}'' .

It can happen that

§18.6 Differentiation of Integrals

Exercises

Bibliography

- [Apo57] T. M. Apostol. *Mathematical Analysis*. Addison-Wesley, 1957.
- [Ax124] S. Axler. *Linear Algebra Done Right*. Springer, 2024.
- [DF04] D. S. Dummit and R. M. Foote. *Abstract Algebra*. John Wiley & Sons, 2004.
- [HS65] E. Hewitt and K. Stromberg. *Real and Abstract Analysis*. Springer-Verlag, 1965.
- [Mun18] J. R. Munkres. *Topology*. Pearson Education Limited, 2018.
- [Pó145] G. Pólya. *How to Solve It*. Princeton University Press, 1945.
- [Rud76] W. Rudin. *Principles of Mathematical Analysis*. McGraw–Hill, 1976.
- [Rud87] W. Rudin. *Real and Complex Analysis*. McGraw–Hill, 1987.
- [Rud91] W. Rudin. *Functional Analysis*. McGraw–Hill, 1991.
- [Sch92] A. H. Schoenfeld. “Learning to think mathematically: Problem solving, metacognition, and sense-making in mathematics”. In: *Handbook for Research on Mathematics Teaching and Learning*. Macmillan, 1992, pp. 334–370.
- [Spi65] M. Spivak. *Calculus on Manifolds: A Modern Approach to Classical Theorems of Advanced Calculus*. Harper Collins Publishers, 1965.

Index

- adjoint, 176
- analytic function, 359
- annihilator, 128
- balls, 214
 - closed ball, 214
 - open ball, 214
 - punctured ball, 214
- basis, 84
- beta functions, 379
- boundary, 219
- boundary point, 219
- boundedness, 214
- Cauchy sequence, 255
- closed set, 217
- closure, 219
- compact, 224
 - open cover, 224
- connectedness, 238
- continuity, 282
 - uniform continuity, 290
- convergence, 246
- coset, 57, 121
 - left coset, 57
 - right coset, 57
- Dedekind cut, 192
- dense, 219
- diagonal matrix, 155
- dimension, 87
- direct sum, 77
- dual basis, 125
- dual map, 126
- eigenspace, 155
- eigenvalue, 141
- eigenvector, 141
- equivalence relation, 29
 - equivalence class, 29
 - partition, 29
 - quotient set, 30
- extended real number system, 202
- finite-dimensional, 81
- Fourier coefficients, 371
- Fourier series, 371
- function, 33
 - bijectivity, 33
 - image, 34
 - injectivity, 33
 - invertibility, 38
 - pre-image, 34
 - restriction, 33
 - surjectivity, 33
- Gamma function, 378
- group, 49
- homomorphism, 63
- image, 64, 96
- index, 58
- induced set, 221
- infimum, 185
- injectivity, 96
- inner product, 160
 - inner product space, 160
- interior, 219
- invariant subspace, 141
- invertibility, 109
- isomorphism, 63, 111
- kernel, 64, 96
- limit of function, 280
- limit point, 221
- linear combination, 80
- linear functional, 125
- linear independence, 81
- linear map, 93
- matrix, 101
 - identity matrix, 115
 - transpose, 107
- matrix of linear map, 101
- matrix of vector, 113
- metric space, 212
- minimal polynomial, 148

- neighbourhood, 216
- normal operator, 178
- normed space, 212

- open set, 216
- operator, 141
- order, 185
- orthogonal projection, 172

- perfect set, 235
- pointwise convergence, 337
- polynomial, 134
 - degree, 134
 - zero, 135
- power series, 358
- product of vector spaces, 118
- pseudoinverse, 174

- quotient map, 122
- quotient space, 121

- rank, 108
 - column rank, 107
 - column space, 107
 - row rank, 107
 - row space, 107
- relation, 28
 - binary relation, 28
 - partial order, 29
 - total order, 29
 - well order, 29
- Riemann–Stieltjes integrability, 317

- self-adjoint operator, 177
- set, 23
 - Cartesian product, 25
 - complement, 26
 - disjoint, 26
 - element, 23
 - empty set, 23
 - intersection, 26
 - interval, 24
 - ordered pair, 25
 - power set, 24
 - set difference, 26
 - subset, 24
 - union, 25
- span, 80
- subgroup, 53
- subsequence, 253
- supremum, 185
- surjectivity, 97

- uniform convergence, 339
- upper-triangular matrix, 152

- vector space, 73
 - complex vector space, 73
 - real vector space, 73
 - subspace, 76