

TOPICS IN
PURE MATHEMATICS

Ryan Joo

The mathematician does not study mathematics because it is useful; he studies it because he delights in it and he delights in it because it is beautiful.

— Henri Poincaré (1854–1912)
French mathematician and theoretical physicist

Copyright © 2025 by Ryan Joo.

This book is licensed under the terms of the Creative Commons Attribution-NonCommercial 4.0 International License (<https://creativecommons.org/licenses/by-nc/4.0>), which permits any noncommercial use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to original author and source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this book are included in the book's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the book's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

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>.

Typeset using L^AT_EX.

Last updated February 24, 2025.

Preface

The reader is not assumed to have any mathematical prerequisites, although some experience with proofs may be helpful. **Preliminary topics** such as logic and methods of proofs (Chapter 1), and basic set theory (Chapter 2) are covered in Part I.

Part II covers **abstract algebra**, which follows [DF04; Art11].

- Chapter 3 introduces groups.

Part III covers **linear algebra**, which follows [Ax124].

- Chapter 4 introduces vector spaces, subspaces, span, linear independence, bases and dimension.
- Chapter 5 concerns linear maps and related concepts.

Part IV covers **real analysis**, which follows [Rud76; Apo57].

- Chapter 6 introduces the real and complex number systems.
- Chapter 7 covers basic point-set topology, in the context of metric spaces.
- Chapter 8 concerns numerical sequences and series, in particular their convergence.
- Chapter 9 covers continuity of functions.
- Chapter 10 covers differentiation.
- Chapter 11 covers Riemann–Stieljes integration.
- Chapter 12 covers sequences and series of functions.
- Chapter 13 covers some special functions, most notably power series and the fourier series.

Part V covers **general topology**, which follows [Mun18].

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 [Pó145], 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 | • Use a model |
| • Look for a pattern | • Consider special cases |
| • Make an orderly list | • Work backwards |
| • Draw a picture | • Use direct reasoning |
| • Eliminate possibilities | • Use a formula |
| • Solve a simpler problem | • Solve an equation |
| • Use symmetry | • 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
1.2	First-order Logic	4
1.3	Methods of Proof	5
	Exercises	14
2	Set Theory	20
2.1	Basics	20
2.2	Relations	25
2.3	Functions	31
2.4	Cardinality	39
II	Abstract Algebra	47
3	Groups	48
3.1	Definition and Properties	48
3.2	Cosets	56
3.3	Homomorphisms and Isomorphisms	62
3.4	Group Actions	66
3.5	Group Product, Finite Abelian Groups	69
	Exercises	70
III	Linear Algebra	71
4	Vector Spaces	72
4.1	Definition of Vector Space	72
4.2	Subspaces	76
4.3	Span and Linear Independence	80
4.4	Bases	84
4.5	Dimension	87

5	Linear Maps	93
5.1	Vector Space of Linear Maps	93
5.2	Kernel and Image	96
5.3	Matrices	101
5.4	Invertibility and Isomorphism	108
5.5	Products and Quotients of Vector Spaces	116
5.6	Duality	121
IV	Real Analysis	126
6	Real and Complex Number Systems	127
6.1	Ordered Sets and Boundedness	127
6.2	Real Numbers	132
6.3	Complex Field	143
6.4	Euclidean Space	147
	Exercises	149
7	Basic Topology	151
7.1	Metric Spaces	152
7.2	Compactness	163
7.3	Perfect Sets	173
7.4	Connectedness	176
	Exercises	178
8	Numerical Sequences and Series	180
8.1	Sequences	180
8.2	Series	195
	Exercises	208
9	Continuity	211
9.1	Limit of Functions	211
9.2	Continuous Functions	214
9.3	Uniform Continuity	221
9.4	Discontinuities	224
9.5	Monotonic Functions	226
9.6	Lipschitz Continuity	228
	Exercises	230

10 Differentiation	232
10.1 The Derivative of A Real Function	232
10.2 Mean Value Theorems	236
10.3 Continuity of Derivatives	238
10.4 L'Hopital's Rule	239
10.5 Taylor's Theorem	240
10.6 Differentiation of Vector-valued Functions	242
Exercises	243
11 Riemann–Stieltjes Integral	244
11.1 Definition of Riemann–Stieltjes Integral	244
11.2 Properties of the Integral	252
11.3 Integration and Differentiation	259
11.4 Integration of Vector-valued Functions	261
11.5 Rectifiable Curves	262
Exercises	264
12 Sequences and Series of Functions	265
12.1 Pointwise Convergence	265
12.2 Uniform Convergence	267
12.3 Properties of Uniform Convergence	270
12.4 Equicontinuous Families of Functions	278
12.5 Stone–Weierstrass Approximation Theorem	282
Exercises	286
13 Some Special Functions	287
13.1 Power Series	287
13.2 Algebraic Completeness of the Complex Field	294
13.3 Fourier Series	295
13.4 Gamma Function	298
Exercises	303
V General Topology	304
14 Topological Spaces and Continuous Functions	305
14.1 Topological Spaces	305
14.2 Basis for a Topology	306

14.3 Examples of Topologies	310
14.4 Closed Sets and Limit Points	313

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, which include

- (i) **equivalence**, denoted by $P \iff Q$, which means P and Q are logically equivalent statements;
- (ii) **conjunction**, denoted by $P \wedge Q$, which means “ P and Q ”;
- (iii) **disjunction**, denoted by $P \vee Q$, which means “ P or Q ”;
- (iv) **negation**, denoted by $\neg P$, which means “not P ”.

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

Proposition 1.1.

- (i) *Double negation law*: $P \iff \neg(\neg P)$.
- (ii) *Commutative*: $P \wedge Q \iff Q \wedge P$, $P \vee Q \iff Q \vee P$.
- (iii) *Conjunction is associative*: $(P \wedge Q) \wedge R \iff P \wedge (Q \wedge R)$.
- (iv) *Disjunction is associative*: $(P \vee Q) \vee R \iff P \vee (Q \vee R)$.
- (v) *Conjunction distributes over disjunction*: $P \wedge (Q \vee R) \iff (P \wedge Q) \vee (P \wedge R)$.
- (vi) *Disjunction distributes over conjunction*: $P \vee (Q \wedge R) \iff (P \vee Q) \wedge (P \vee R)$.

Proposition 1.2 (de Morgan's laws).

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

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

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* 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 $P \implies Q$, start by assuming that P holds and try to deduce through some logical steps that Q holds too. Alternatively, start by assuming that Q does not hold and show that P does not hold (that is, we prove the contrapositive).

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, i.e. prove both $P \implies Q$ and $Q \implies P$. 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*.

Proposition 1.3 (de Morgan’s laws).

$$\neg \forall x \in X, P(x) \iff \exists x \in X, \neg P(x)$$

$$\neg \exists x \in X, P(x) \iff \forall x \in X, \neg P(x)$$

To prove a statement of the form $\forall x \in X$ s.t. $P(x)$, start the proof with “Let $x \in X$.” or “Suppose $x \in X$ is given.” to address the quantifier with an arbitrary x ; provided no other assumptions about x are made during the course of proving $P(x)$, this will prove the statement for all $x \in X$.

To prove a statement of the form $\exists x \in X$ s.t. $P(x)$, there is not such a clear steer about how to continue: you may need to show the existence of an x with the right properties; you may need to demonstrate logically that such an x must exist because of some earlier assumption, or it may be that you can show constructively how to find one; or you may be able to prove by contradiction, supposing that there is no such x and consequently arriving at some inconsistency.

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

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 . It may be that we can start from P and work directly to Q , or it may be that we make use of P along the way.

A **proof by contrapositive** of $P \implies Q$ is to prove instead $\neg Q \implies \neg P$.

A **disproof by counterexample** is to providing a counterexample in order to refute or disprove a conjecture. The counterexample must make the hypothesis a true statement, and the conclusion a false statement. 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

A **proof by contradiction** of P involves first supposing P is false, i.e. $\neg P$ (to prove $P \implies Q$ by contradiction, suppose $P \wedge \neg Q$). Then 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.

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

To prove existential statements, we can adopt two approaches:

1. **Constructive proof** (direct proof)

To prove statements of the form $\exists x \in X$ s.t. $P(x)$, find or construct a *specific example* for x . To prove statements of the form $\forall y \in Y, \exists x \in X$ s.t. $P(x, y)$, 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**, we can either assume $\exists x, y \in S$ such that $P(x) \wedge P(y)$ is true and show $x = y$, or argue by assuming that $\exists x, y \in S$ are distinct such that $P(x) \wedge P(y)$, then derive a contradiction. $\exists!$ denotes “there exists a unique”. To prove uniqueness and existence, we also need to show that $\exists x \in S$ s.t. $P(x)$ is true.

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 $r = p + \frac{q-p}{\sqrt{2}} > p + 0 = p$.

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

(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 $(a, b, c) = \left(\frac{x}{3}, \frac{y}{3}, \frac{z}{3}\right)$, 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 classical 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

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

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 $N \in \mathbb{Z}$. 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 $k \in \mathbb{N}$, $P(1) \wedge \dots \wedge P(k) \implies P(k + 1)$.

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 $P(k) \implies (k-1)$.

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 ,

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

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

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 $P(n) \implies P(2n)$

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

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, $\frac{a_1 + a_2 + \dots + a_n}{n} \geq \sqrt[n]{a_1 a_2 \dots a_n}$

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. \square

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

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

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 $x_1 + x_2 + \dots + x_k$ and $x_1 + x_2 + \dots + x_{k+m}$, belong to the same interval. But then their difference $x_{k+1} + \dots + x_{k+m}$ lies within a distance of $\frac{1}{n+1}$ of an integer, and we are done.

□

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 $(\forall x \in \mathbb{R})[x > 2 \implies x^2 > 4]$.

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

□

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 $f(x) \neq y$.

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

□

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 . Evidently $p < \frac{p+q}{2} < q$.

- If $\frac{p+q}{2}$ is irrational, take $x = \frac{p+q}{2}$ and we are done.
- If $\frac{p+q}{2}$ is rational, call it r , take the average of p and r : $p < \frac{p+r}{2} < r < q$. Since p is irrational and r is rational, $\frac{p+r}{2}$ is irrational. In this case, we 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 $\frac{\sqrt{2}}{a}$.

Case 1: $\frac{\sqrt{2}}{a}$ is irrational.

Take $b = \frac{\sqrt{2}}{a}$. Then $ab = \sqrt{2}$ which is irrational.

Case 2: $\frac{\sqrt{2}}{a}$ is rational.

Then the reciprocal $\frac{a}{\sqrt{2}}$. 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 $ab = a\sqrt{3}$ which is irrational. \square

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

Solution. Prove by contradiction.

Suppose 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.

We construct an integer M by $M = (p_1 p_2 \cdots p_m)^2 + 2$.

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. \square

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

Solution. Prove by contradiction.

Suppose otherwise, that $n > m^2$ and $(m+1)^2 > 2n$ 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. \square

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

$$n! > 2^n.$$

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

For the base case $P(4)$, the LHS equals to $4! = 4 \times 3 \times 2 \times 1 = 24$, and the RHS equals $2^4 = 16 < 24$. Since LHS equals RHS, $P(4)$ is true.

Now suppose that $P(k)$ is true for some $k \in \mathbb{N}_{\geq 4}$. Then

$$\begin{aligned} k! &> 2^k \\ (k+1)k! &> 2^k(k+1) \\ &> 2^k 2 \quad \text{since from } k \geq 4, k+1 \geq 5 > 2 \\ &= 2^{k+1} \end{aligned}$$

thus $P(k+1)$ is true, so we have shown $P(k) \implies P(k+1)$ for all $k \in \mathbb{N}_{\geq 4}$.

By PMI, we have proven $P(n)$ for all integers $n \geq 4$. □

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

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

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

For the base case, when $n = 2$, $\sqrt{2} = 1.41 \dots < 2 - \frac{1}{2} = 1.5$ which is true. Hence $P(2)$ is true.

Now assume $P(k)$ is true for $k \geq 2, k \in \mathbb{N}$; that is,

$$\sqrt[k]{k} < 2 - \frac{1}{k} \implies k < \left(2 - \frac{1}{k}\right)^k$$

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

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

Hence $P(k+1)$ is true.

Since $P(2)$ is true and $P(k) \implies P(k+1)$, by mathematical induction $P(n)$ 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. □

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

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

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

Given the recurrence relation, it could be possible to use $P(k)$, $P(k+1)$, $P(k+2)$ to prove $P(k+3)$ for all $k \in \mathbb{N}$.

Base case: $P(1), P(2), P(3)$

$P(1) : a_1 = 1 \leq 2^{1-1} = 1$ is true.

$P(2) : a_2 = 1 \leq 2^{2-1} = 2$ is true.

$P(3) : a_3 = 3 \leq 2^{3-1} = 4$ is true.

Inductive step: $P(k) \wedge P(k+1) \wedge P(k+2) \implies P(k+3)$ for all $k \in \mathbb{N}$

By inductive hypothesis, for $k \in \mathbb{N}$ we have $a_k \leq 2^k, a_{k+1} \leq 2^{k+1}, a_{k+2} \leq 2^{k+2}$.

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

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

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. We 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

$$F_{m+1} = F_0 F_m + F_1 F_{m+1}$$

and

$$F_{m+2} = F_1 F_m + F_2 F_{m+1}$$

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

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

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

Then

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

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

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:

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

Some basic sets (of numbers) you should be familiar with:

- $\mathbb{N} = \{1, 2, 3, \dots\}$ 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 6).
- $\mathbb{C} = \{x + yi \mid x, y \in \mathbb{R}\}$ 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 .

¹Russell'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 , denoted by $A \times B$, 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.

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

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

- (i) $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$;
- (ii) $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$.

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

- (i) $(A \cup B)^c = A^c \cap B^c$;
- (ii) $(A \cap B)^c = A^c \cup B^c$.

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

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

Exercise 2.1. Prove the following:

- (i) $\left(\bigcup_{i \in I} A_i\right) \cup B = \bigcup_{i \in I} (A_i \cup B)$
- (ii) $\left(\bigcap_{i \in I} A_i\right) \cup B = \bigcap_{i \in I} (A_i \cup B)$
- (iii) $\left(\bigcup_{i \in I} A_i\right) \cup \left(\bigcup_{j \in J} B_j\right) = \bigcup_{(i,j) \in I \times J} (A_i \cup B_j)$
- (iv) $\left(\bigcap_{i \in I} A_i\right) \cup \left(\bigcap_{j \in J} B_j\right) = \bigcap_{(i,j) \in I \times J} (A_i \cup B_j)$

Exercise 2.2. Let $S \subset A \times B$. Express the set A_S of all elements of A which appear as the first entry in at least one of the elements in S .

(A_S here may be called the projection of S onto A .)

§2.2 Relations

Definition and Examples

Definition 2.6 (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 2.7. 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 \mid 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 2.8. Given $R = \{(1, a), (1, b), (2, b), (3, b)\}$, then $\text{dom } R = \{1, 2, 3\}$ and $\text{ran } R = \{a, b\}$.

Definition 2.9 (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 2.10 (Less than or equal to). The relation \leq on R is reflexive, anti-symmetric, and transitive, but not symmetric.

Definition 2.11. 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 2.12.

- 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 2.13 (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 2.14 (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 2.15 (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

- (i) $A_i \neq \emptyset$ for all $i \in I$ (all the subsets are non-empty)

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

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

Proposition 2.16. *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 2.17 (Quotient set). The *quotient set* is the set of all equivalence classes, denoted by A/\sim .

Example 2.18 (Modular arithmetic). Let n be a fixed positive integer. Define a relation on \mathbb{Z} by

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

Proposition. $a \sim b$ is a equivalence relation.

Proof.

- (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 $a \sim b$.

Notation. For any $k \in \mathbb{Z}$ we denote the equivalence class of a by $[a]$, called the *congruence class* (or *residue class*) of $a \pmod{n}$, 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 $[a], [b] \in \mathbb{Z}/n\mathbb{Z}$,

$$[a] + [b] = [a + b]$$

$$[a][b] = [ab].$$

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.

Proposition. *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. □

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

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

Proposition. $(\mathbb{Z}/n\mathbb{Z})^\times$ is also 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 2.19. 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 2.20. 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 2.21. 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 2.22. 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. □

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

Definition 2.23 (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 2.24 (Restriction). Given a function $f : X \rightarrow Y$ and a subset $A \subset X$, the **restriction** of f to A 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 2.25 (Identity map). Given a set X , the **identity** $\text{id}_X : X \rightarrow X$ is defined by

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

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

Images and Pre-images

Definition 2.26. Suppose $f : X \rightarrow Y$. The **image** of f is

$$f(X) := \{f(x) \mid x \in X\} \subset Y.$$

More generally, the image of $A \subset X$ under f is

$$f(A) := \{f(x) \mid x \in A\} \subset Y.$$

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 2.27. 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 2.28 (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 $f^{-1} \left(\bigcap_{i \in I} A_i \right) = \bigcap_{i \in I} f^{-1}(A_i)$.

□

Injectivity, Surjectivity, Bijectivity

Definition 2.29. 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 $f : X \rightarrow Y$.

Composition

Definition 2.30 (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 2.31 (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 2.32 (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 2.33. $f : X \rightarrow Y$ 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).

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

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

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 2.34. $f : X \rightarrow Y$ 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.

\implies 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 2.35 (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 2.36 (Invertibility). Suppose $f : X \rightarrow Y$. We say that

- (i) f is **left-invertible** if there exists $g : Y \rightarrow X$ such that $g \circ f = \text{id}_X$; g is a *left-inverse* of f ;
- (ii) f is **right-invertible** if there exists $h : Y \rightarrow X$ such that $f \circ h = \text{id}_Y$; h is 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 ; k is 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 2.37. 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

Proposition 2.38 (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. Note that directly 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 2.39 (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 2.40 (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 2.41. *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 2.42. *\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 2.43 (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 2.44. 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

$$\mathbb{N}_n = \{k \in \mathbb{N} \mid 1 \leq k \leq n\},$$

$$n\mathbb{N} = \{nk \mid k \in \mathbb{N}\}.$$

Definition 2.45. 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.

Proposition 2.46. \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 2.47. Any countable set is infinite.

Example 2.48. \mathbb{N} is countable since the identity map from \mathbb{N} to \mathbb{N} is a bijection.

Example 2.49. $n\mathbb{N}$ is countable.

Proof. Let $f : \mathbb{N} \rightarrow n\mathbb{N}$ which sends $k \mapsto nk$. We now need to show that f is (i) injective, and (ii) surjective.

- (i) For any $k_1, k_2 \in \mathbb{N}$, $nk_1 = nk_2$ implies $k_1 = k_2$ so f is injective.

(ii) 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. \square

Example 2.50. \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

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

\square

Proposition 2.51. *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. \square

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

Proposition 2.52. *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.

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

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 2.51, 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 2.53. 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 2.54. 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 2.52, B_n is countable. \square

Corollary 2.55. \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 2.56. 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 2.57. \mathbb{R} is uncountable.

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

Theorem 2.58 (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 $A \rightarrow \mathcal{P}(A)$. □

Exercises

Exercise 2.3. Prove the following statements:

- (i) $f(A \cup B) = f(A) \cup f(B)$
- (ii) $f(A_1 \cup \dots \cup A_n) = f(A_1) \cup \dots \cup f(A_n)$
- (iii) $f(\bigcup_{\lambda \in A} A_\lambda) = \bigcup_{\lambda \in A} f(A_\lambda)$
- (iv) $f(A \cap B) \subset f(A) \cap f(B)$
- (v) $f^{-1}(f(A)) \supset A$
- (vi) $f(f^{-1}(A)) \subset A$
- (vii) $f^{-1}(A \cup B) = f^{-1}(A) \cup f^{-1}(B)$
- (viii) $f^{-1}(A \cap B) = f^{-1}(A) \cap f^{-1}(B)$
- (ix) $f^{-1}(A_1 \cup \dots \cup A_n) = f^{-1}(A_1) \cup \dots \cup f^{-1}(A_n)$
- (x) $f^{-1}(\bigcup_{\lambda \in A} A_\lambda) = \bigcup_{\lambda \in A} f^{-1}(A_\lambda)$

Exercise 2.4. 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

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

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

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

Proof. 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 $g \in T$, 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 2.5. 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)

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

$$\text{Thus } [x' + y'] = [x + y]$$

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

$$\text{Thus } [-x'] = [-x].$$

- (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.

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

\square

Exercise 2.6 (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 2.7 ([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 2.8 ([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 2.9 ([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) Associativity: $a(bc) = (ab)c$ for all $a, b, c \in G$;
- (ii) Identity: there exists $e \in G$ such that $ae = ea = a$ for all $a \in G$.
- (iii) Invertibility: for all $a \in G$, there exists $c \in G$ such that $ac = ca = e$.

G is *abelian* if the operation is commutative; it is *non-abelian* if otherwise.

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 and under multiplication as $(\mathbb{Z}_n)^\times$.

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 is an abelian group under $+$.
- For $n \in \mathbb{Z}^+$, $(\mathbb{Z}_n)^\times$ is an abelian group under multiplication.

Proposition 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.

□

Proposition 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 $ya = b$. □

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 $rsr^{-k} = sr^{-1}r^{-k} = sr^{-k-1}$ 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.

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

\Leftarrow 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. □

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.

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

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

⊇ 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 $o(a) \mid n$. □

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.

\Rightarrow Suppose $aH = H$. Then $ah \in H$ for some $h \in H$. Let $k = ah$, then $a = kh^{-1} \in H$.

\Leftarrow 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 $G = \bigcup_{a \in G} aH$. 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 $|G|$. \square

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 $x^{-1} \in HK$.

\Rightarrow 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 $gxg^{-1} = h'gg^{-1} = h' \in H$.

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

$$\operatorname{im} \phi := \phi(G) = \{\phi(g) \mid g \in G\} \subset H.$$

Remark. $\operatorname{im} \phi$ is the usual set theoretic image of ϕ .

Proposition 3.52. *Let $\phi : G \rightarrow H$ be a homomorphism. Then*

$$(i) \quad \ker \phi \triangleleft G;$$

$$(ii) \quad \operatorname{im} \phi \leq H.$$

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 $\phi(1_G) = 1_H$, $1_H \in \operatorname{im} \phi$ so $\operatorname{im} \phi \neq \emptyset$. Let $x, y \in \operatorname{im} \phi$. 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 $xy^{-1} \in \operatorname{im} \phi$. By the subgroup criterion, $\operatorname{im} \phi \leq G$.

□

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 $\text{Syl}_p(G)$, 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. $\text{Syl}_p(G) \neq \emptyset$.
- (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 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*.

Proposition 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

Proposition 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.

Proposition 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) For $v \in V$, we have

$$0v = (0 + 0)v = 0v + 0v.$$

Adding the additive inverse of $0v$ to both sides of the equation gives $\mathbf{0} = 0v$.

- (ii) For $a \in \mathbf{F}$, we have

$$a\mathbf{0} = a(\mathbf{0} + \mathbf{0}) = a\mathbf{0} + a\mathbf{0}.$$

Adding the additive inverse of $a\mathbf{0}$ to both sides of the equation gives $\mathbf{0} = a\mathbf{0}$.

- (iii) For $v \in V$, we have

$$v + (-1)v = 1v + (-1)v = (1 + (-1))v = 0v = \mathbf{0}.$$

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} :

$$\mathbf{F}^\infty := \{(x_1, x_2, \dots) \mid x_i \in \mathbf{F}\}$$

- Addition on \mathbf{F}^∞ is defined by

$$(x_1, x_2, \dots) + (y_1, y_2, \dots) = (x_1 + y_1, x_2 + y_2, \dots)$$

- Scalar multiplication on \mathbf{F}^∞ is defined by

$$\lambda(x_1, x_2, \dots) = (\lambda x_1, \lambda x_2, \dots)$$

Verify that \mathbf{F}^∞ becomes a vector space over \mathbf{F} . Also verify that the additive identity in \mathbf{F}^∞ 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, $\mathbf{F}^S := \{f \mid f : S \rightarrow \mathbf{F}\}$.

- Addition on \mathbf{F}^S is defined by

$$(f + g)(x) = f(x) + g(x) \quad (\forall x \in S)$$

for all $f, g \in \mathbf{F}^S$.

- Multiplication on \mathbf{F}^S is defined by

$$(\lambda f)(x) = \lambda f(x) \quad (\forall x \in S)$$

for all $\lambda \in \mathbf{F}, f \in \mathbf{F}^S$.

Verify that if S is a non-empty set, then \mathbf{F}^S is a vector space over \mathbf{F} .

Also verify that the additive identity of \mathbf{F}^S is the function $0 : S \rightarrow \mathbf{F}$ defined by

$$0(x) = 0 \quad (\forall x \in S)$$

and for $f \in \mathbf{F}^S$, additive inverse of f is the function $-f : S \rightarrow \mathbf{F}$ defined by

$$(-f)(x) = -f(x) \quad (\forall x \in S)$$

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 by

$$(u_1 + iv_1) + (u_2 + iv_2) = (u_1 + u_2) + i(v_1 + v_2)$$

for all $u_1, v_1, u_2, v_2 \in V$.

- Complex scalar multiplication on $V_{\mathbb{C}}$ is defined by

$$(a + bi)(u + iv) = (au - bv) + i(av + bu)$$

for all $a, b \in \mathbb{R}$ and all $u, v \in V$.

You should verify that with the definitions of addition and scalar multiplication as above, $V_{\mathbb{C}}$ 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). We denote this 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) *Additive identity*: $0 \in U$
- (ii) *Closed under addition*: $u + w \in U$ for all $u, w \in U$
- (iii) *Closed under scalar multiplication*: $\lambda u \in U$ for all $\lambda \in \mathbf{F}, u \in U$

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 . \square

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 $U = \{(x, 0, 0) \in \mathbf{F}^3 \mid x \in F\}$ 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 $v_1 = \dots = v_{n-1} = v_n = \mathbf{0}$. □

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 $\mathbf{0} = u + w$. $u = -w \in W$, thus $u \in U \cap W$, so $u = w = \mathbf{0}$. By Lemma 4.17, $U + W$ is a direct sum. □

§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

$$v = a_1 v_1 + \dots + a_n v_n.$$

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 :

$$\text{span}(v_1, \dots, v_n) := \{a_1 v_1 + \dots + a_n v_n \mid a_i \in \mathbf{F}\}.$$

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 $\text{span}(S) = V$, then we say that S *spans* V , and that S is a *spanning set* for V :

$$\text{span}(S) := \{a_1 v_1 + \dots + a_n v_n \mid v_i \in S, a_i \in \mathbf{F}\}.$$

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) $(a_1 v_1 + \dots + a_n v_n) + (c_1 v_1 + \dots + c_n v_n) = (a_1 + c_1)v_1 + \dots + (a_n + c_n)v_n \in \text{span}(v_1, \dots, v_n)$, so $\text{span}(v_1, \dots, v_n)$ is closed under addition.
- (iii) $\lambda(a_1 v_1 + \dots + a_n v_n) = (\lambda a_1)v_1 + \dots + (\lambda a_n)v_n \in \text{span}(v_1, \dots, v_n)$, 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

$$a_1 v_1 + \dots + a_n v_n \in M$$

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 vector space). 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.

Proposition 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:

$$(i) \quad v_k \in \text{span}(v_1, \dots, v_{k-1})$$

$$(ii) \quad \text{span}(v_1, \dots, v_{k-1}, v_{k+1}, \dots, v_n) = \text{span}(v_1, \dots, v_n)$$

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

$$a_1 v_1 + \dots + a_n v_n = 0.$$

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

$$v_k = b_1 v_1 + \dots + b_{k-1} v_{k-1}. \quad (1)$$

Suppose $u \in \text{span}(v_1, \dots, v_n)$. Then there exists $c_1, \dots, c_n \in \mathbf{F}$ such that

$$u = c_1 v_1 + \dots + c_n v_n. \quad (2)$$

In (2), we can replace v_k with the RHS of (1), which gives

$$\begin{aligned} u &= c_1 v_1 + \dots + c_{k-1} v_{k-1} + c_k v_k + c_{k+1} v_{k+1} + \dots + c_n v_n \\ &= c_1 v_1 + \dots + c_{k-1} v_{k-1} + c_k (b_1 v_1 + \dots + b_{k-1} v_{k-1}) + c_{k+1} v_{k+1} + \dots + c_n v_n \\ &= c_1 v_1 + \dots + c_{k-1} v_{k-1} + c_k b_1 v_1 + \dots + c_k b_{k-1} v_{k-1} + c_{k+1} v_{k+1} + \dots + c_n v_n \\ &= (c_1 + b c_k) v_1 + \dots + (c_{k-1} + b_{k-1} c_k) v_{k-1} + c_{k+1} v_{k+1} + \dots + c_n v_n. \end{aligned}$$

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 $U = \{0\}$, then U is finite-dimensional and we are done. Otherwise, choose $v_1 \in U$, $v_1 \neq 0$ 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 Proposition 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 $\mathbf{e}_i = (0, \dots, 0, 1, 0, \dots, 0)$ where the i -th coordinate is 1. $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$ 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 list 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

$$v = a_1 v_1 + \dots + a_n v_n. \quad (1)$$

To show that the representation is unique, suppose that $c_1, \dots, c_n \in \mathbf{F}$ also satisfy

$$v = c_1 v_1 + \dots + c_n v_n. \quad (2)$$

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 ($1 \leq i \leq n$). 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 $a_1, \dots, a_n \in \mathbf{F}$ satisfy

$$a_1 v_1 + \dots + a_n v_n = \mathbf{0}.$$

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. Our 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 Lemma 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 Lemma 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 Proposition 4.27, U is finite-dimensional, so U has a basis B , by Corollary 4.32; let $B = \{u_1, \dots, u_n\}$. Since B is linearly independent, by Lemma 4.33, B can be extended to a basis of V , say

$$\{u_1, \dots, u_n, w_1, \dots, w_n\}.$$

Take $W = \text{span}(w_1, \dots, w_n)$. We claim that $V = U \oplus W$. To show this, by Lemma 4.17, we need to show that (i) $V = U + W$, and (ii) $U \cap W = \{\mathbf{0}\}$.

(i) Suppose $v \in V$. 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) Suppose $v \in U \cap W$. 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 Proposition 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 Proposition 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 Lemma 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 Lemma 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.

□

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 Proposition 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 $U = V$. □

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 Lemma 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, we have $a_i = b_i = 0$ for all i , as desired. □

Exercises

Exercise 4.2 ([Axl24] 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 ([Axl24] 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 ([Axl24] 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,

$$a_1v_1 + a_2(v_1 + v_2) + \cdots + a_n(v_1 + \cdots + v_n) = 0,$$

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 Proposition 4.26, we have reached a contradiction. \square

Exercise 4.7 ([Ax124] 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 ([Ax124] 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

$$a_1(v_1 + v_2) + a_2(v_2 + v_3) + a_3(v_3 + v_4) + a_4 v_4 = \mathbf{0} \implies a_1 = a_2 = a_3 = a_4 = 0.$$

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

$$\begin{aligned} v &= a_1v_1 + a_2v_2 + a_3v_3 + a_4v_4 \\ &= a_1(v_1 + v_2) + (a_2 - a_1)(v_2 + v_3) + (a_3 - a_2)(v_3 + v_4) + (a_4 - a_3)v_4, \end{aligned}$$

which is a linear combination of $v_1 + v_2, v_2 + v_3, v_3 + v_4, v_4$. □

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) Additivity: $T(v + w) = Tv + Tw$ for all $v, w \in V$
- (ii) Homogeneity: $T(\lambda v) = \lambda T(v)$ for all $\lambda \in \mathbf{F}, v \in V$

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 Lemma 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 $c_1v_1 + \dots + c_nv_n$,

$$\begin{aligned} T(u + v) &= T((a_1 + c_1)v_1 + \dots + (a_n + c_n)v_n) \\ &= (a_1 + c_1)w_1 + \dots + (a_n + c_n)w_n \\ &= (a_1w_1 + \dots + a_nw_n) + (c_1w_1 + \dots + c_nw_n) \\ &= Tu + Tv. \end{aligned}$$

(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 $S, T \in \mathcal{L}(V, W)$, $\lambda \in \mathbf{F}$,

$$(i) \quad (S + T)(v) = Sv + Tv$$

$$(ii) \quad (\lambda T)(v) = \lambda(Tv)$$

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

$$(ST)(u) = S(Tu) \quad (\forall u \in U)$$

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 Lemma 4.10, we check the conditions of a subspace:

(i) By Proposition 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 $\operatorname{im} T$ is finite-dimensional, and*

$$\dim V = \dim \ker T + \dim \operatorname{im} T. \quad (5.1)$$

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 $\operatorname{im} T$ is finite-dimensional and $\dim \operatorname{im} T = n$, we prove that $\{Tv_1, \dots, Tv_n\}$ is a basis of $\operatorname{im} T$. Thus we need to show that the set (i) spans $\operatorname{im} T$, 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

$$v = a_1 u_1 + \dots + a_m u_m + b_1 v_1 + \dots + b_n v_n,$$

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

$$u_i \in \ker T,$$

$$\begin{aligned} Tv &= T(a_1u_1 + \cdots + a_mu_m + b_1v_1 + \cdots + b_nv_n) \\ &= a_1 \underbrace{Tu_1}_{\mathbf{0}} + \cdots + a_m \underbrace{Tu_m}_{\mathbf{0}} + b_1Tv_1 + \cdots + b_nv_n \\ &= b_1Tv_1 + \cdots + b_nv_n \in \operatorname{im} T. \end{aligned}$$

Since every element of $\operatorname{im} T$ can be expressed as a linear combination of Tv_1, \dots, Tv_n , we have that $\{Tv_1, \dots, Tv_n\}$ spans $\operatorname{im} T$.

Moreover, since there exists a set of vectors that spans $\operatorname{im} T$, $\operatorname{im} T$ is finite-dimensional.

(ii) Suppose there exist $c_1, \dots, c_n \in \mathbf{F}$ such that

$$c_1Tv_1 + \cdots + c_nTv_n = \mathbf{0}.$$

Then

$$T(c_1v_1 + \cdots + c_nv_n) = T(\mathbf{0}) = \mathbf{0},$$

which implies $c_1v_1 + \cdots + c_nv_n \in \ker T$. Since $\{u_1, \dots, u_m\}$ is a spanning set of $\ker T$, we can write

$$c_1v_1 + \cdots + c_nv_n = d_1u_1 + \cdots + d_mu_m$$

for some $d_i \in \mathbf{F}$, or

$$c_1v_1 + \cdots + c_nv_n - d_1u_1 - \cdots - d_mu_m = \mathbf{0}.$$

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 $\operatorname{im} T \leq W$, by Proposition 4.37, we have that $\dim \operatorname{im} T \leq \dim W$.

Let $T \in \mathcal{L}(V, W)$. Then

$$\dim \ker T = \dim V - \dim \operatorname{im} T \tag{1}$$

$$\geq \dim V - \dim W \tag{2}$$

$$> 0$$

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 $v \in V \setminus \{\mathbf{0}\}$. Since $\ker T \neq \{\mathbf{0}\}$, T is not injective. □

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

$$\dim \operatorname{im} T = \dim V - \dim \ker T \quad (1)$$

$$\leq \dim V \quad (2)$$

$$< \dim W,$$

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 $\dim \operatorname{im} T < \dim W$, $\operatorname{im} T \neq W$ so T is not surjective. □

Example 5.17 (Homogeneous system of linear equations). Consider the homogeneous system of linear equations

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &= 0 \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n &= 0 \\ &\vdots \\ a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n &= 0 \end{aligned} \quad (*)$$

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

$$\begin{aligned} \dim \mathbf{F}^n > \dim \mathbf{F}^m &\implies T \text{ is not injective} \\ &\implies \ker T \neq \{0\} \\ &\implies (*) \text{ has non-zero solutions} \end{aligned}$$

□

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

$$\dim \mathbf{F}^n < \dim \mathbf{F}^m \implies T \text{ is not surjective}$$

$$\implies \exists (c_1, \dots, c_m) \in \mathbf{F}^m, \forall (x_1, \dots, x_n) \in \mathbf{F}^n, T(x_1, \dots, x_n) \neq (c_1, \dots, c_m)$$

Thus the choice of constant terms (c_1, \dots, c_m) is such that the system of linear equations

$$\begin{aligned} a_{11}x_1 + \cdots + a_{1n}x_n &= c_1 \\ &\vdots \\ a_{m1}x_1 + \cdots + a_{mn}x_n &= c_m \end{aligned}$$

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:

$$\lambda \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix} = \begin{pmatrix} \lambda a_{11} & \cdots & \lambda a_{1n} \\ \vdots & & \vdots \\ \lambda a_{m1} & \cdots & \lambda a_{mn} \end{pmatrix}.$$

Proposition 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

$$\begin{aligned}\lambda T v_j &= \sum_{i=1}^m b_{ij} w_i \\ \lambda \sum_{i=1}^m a_{ij} w_i &= \sum_{i=1}^m b_{ij} w_i \\ \lambda a_{ij} &= b_{ij}.\end{aligned}$$

□

Proposition 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

Definition 5.23 (Matrix multiplication). Suppose $A = (a_{ij})_{m \times n}$, $B = (b_{jk})_{n \times p}$. Then $AB = (c_{ik})_{m \times p}$ has entries defined by

$$c_{ik} = \sum_{j=1}^n a_{ij} b_{jk}.$$

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.

Remark. 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.

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 $ST \in \mathcal{L}(U, W)$.

Proposition 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 $\mathcal{M}(S) = (a_{ij})_{m \times n}$, $\mathcal{M}(T) = (b_{jk})_{n \times p}$, where

$$Sv_j = \sum_{i=1}^m a_{ij}w_i$$

$$Tu_k = \sum_{j=1}^n b_{jk}v_j.$$

For $k = 1, \dots, p$, we have

$$\begin{aligned} (ST)u_k &= S(Tu_k) \\ &= S\left(\sum_{j=1}^n b_{jk}v_j\right) \\ &= \sum_{j=1}^n b_{jk}Sv_j \\ &= \sum_{j=1}^n b_{jk}\left(\sum_{i=1}^m a_{ij}w_i\right) \\ &= \sum_{i=1}^m \left(\sum_{j=1}^n a_{ij}b_{jk}\right)w_i. \end{aligned}$$

□

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 .

Proposition 5.25. Suppose $A = (a_{ij})_{m \times n}$, $B = (b_{jk})_{n \times p}$. Let $AB = (c_{ik})_{m \times p}$. Then

$$c_{ik} = A_{i,\cdot}B_{\cdot,k}$$

That is, the entry in row i , column k of AB equals (row i of A) times (column k of B).

Proof. By definition,

$$\begin{aligned} A_{i,\cdot}B_{\cdot,k} &= \begin{pmatrix} a_{i1} & \cdots & a_{in} \end{pmatrix} \begin{pmatrix} b_{1k} \\ \vdots \\ b_{nk} \end{pmatrix} \\ &= a_{i1}b_{1k} + \cdots + a_{in}b_{nk} \\ &= \sum_{j=1}^n a_{ij}b_{jk} \\ &= c_{ik}. \end{aligned}$$

□

Proposition 5.26. Suppose $A = (a_{ij})_{m \times n}$, $B = (b_{jk})_{n \times p}$. Then

$$(AB)_{\cdot, k} = AB_{\cdot, k}$$

That is, column k of AB equals A times column k of B .

Proof. Using the previous result,

$$AB_{\cdot, k} = \begin{pmatrix} A_{1, \cdot} B_{\cdot, k} \\ \vdots \\ A_{n, \cdot} B_{\cdot, k} \end{pmatrix} = \begin{pmatrix} c_{1k} \\ \vdots \\ c_{nk} \end{pmatrix} = (AB)_{\cdot, k}$$

□

Proposition 5.27 (Linear combination of columns). Suppose $A = (a_{ij})_{m \times n}$, $b = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix}$. Then

$$Ab = b_1 A_{\cdot, 1} + \cdots + b_n A_{\cdot, n}.$$

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 .

Proposition 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

$$S_1 = S_1 I_W = S_1 (TS_2) = (S_1 T) S_2 = I_V S_2 = S_2.$$

Thus $S_1 = S_2$. □

Now that we know that 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, $\dim V = \dim W$, $S \in \mathcal{L}(W, V)$, $T = \mathcal{L}(V, W)$. 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 $\operatorname{im} T = W$, implying $\dim \ker T = 0$ and $\dim \operatorname{im} T = \dim W$.

By the fundamental theorem of linear maps,

$$\begin{aligned} \dim V &= \dim \ker T + \dim \operatorname{im} T \\ &= 0 + \dim W = \dim W. \end{aligned}$$

\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 $w = a_1w_1 + \dots + a_nw_n$. Then

$$\begin{aligned} T(a_1v_1 + \dots + a_nv_n) &= w \implies w \in \operatorname{im} T \\ &\implies W = \operatorname{im} T \\ &\implies T \text{ is surjective} \\ &\implies T \text{ is invertible.} \end{aligned}$$

□

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 $T = 0$, so $\ker \mathcal{M} = \{0\}$. 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 $\mathcal{M}(T) = A$, $\operatorname{im} \mathcal{M} = \mathcal{M}_{m \times n}(\mathbf{F})$ so \mathcal{M} is surjective.

□

Corollary 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

$$\mathcal{M}(x) = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}.$$

Proposition 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

$$\mathcal{M}(T)_{\cdot, j} = \mathcal{M}(Tv_j) \quad (j = 1, \dots, n)$$

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 $\mathcal{M}(T)$, for $j = 1, \dots, n$. □

The following result shows that linear maps act like matrix multiplication.

Proposition 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 $v = b_1v_1 + \dots + b_nv_n$ for some $b_1, \dots, b_n \in \mathbf{F}$. Then

$$\begin{aligned} \mathcal{M}(Tv) &= \mathcal{M}(T(b_1v_1 + \dots + b_nv_n)) \\ &= b_1\mathcal{M}(Tv_1) + \dots + b_n\mathcal{M}(Tv_n) \\ &= b_1\mathcal{M}(T)_{\cdot,1} + \dots + b_n\mathcal{M}(T)_{\cdot,n} \\ &= \begin{pmatrix} \mathcal{M}(T)_{\cdot,1} & \dots & \mathcal{M}(T)_{\cdot,n} \end{pmatrix} \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} \\ &= \mathcal{M}(T)\mathcal{M}(v). \end{aligned}$$

□

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.

Proposition 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

Proposition 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. Suppose otherwise, for a contradiction, that A does not have a unique inverse. Let B and C be inverses of A ; that is,

$$\begin{aligned} AB &= BA = I, \\ AC &= CA = I. \end{aligned}$$

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} .

Proposition 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)

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

so the inverse of A^{-1} is A .

(ii)

$$\begin{aligned} (AC)(C^{-1}A^{-1}) &= A(CC^{-1})A^{-1} \\ &= AIA^{-1} \\ &= AA^{-1} \\ &= I, \end{aligned}$$

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.

Proof. □

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. □

Proposition 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.

Proposition 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 $V_1 + \dots + V_n$ 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.

Proposition 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

Proposition 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. \square

Definition 5.67 (Quotient map). Suppose $U \leq V$. The *quotient map* $\pi : V \rightarrow V/U$ is the linear map defined by

$$\pi(v) = v + U$$

for all $v \in V$.

Proposition 5.68 (Dimension of quotient space). Suppose V is finite-dimensional, $U \leq V$. Then

$$\dim V/U = \dim V - \dim U.$$

Definition 5.69. Suppose $T \in \mathcal{L}(V, W)$. Define $\tilde{T} : V/\ker T \rightarrow W$ by

$$\tilde{T}(v + \ker T) = Tv.$$

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) $\operatorname{im} \tilde{T} = \operatorname{im} T$.

Theorem 5.71 (First isomorphism theorem). *Suppose $T \in \mathcal{L}(V, W)$ is an isomorphism. Then*

$$V/\ker T \cong \operatorname{im} T. \quad (5.3)$$

§5.6 Duality

Dual Space and Dual Map

Linear maps into the scalar field \mathbf{F} play a special role in linear algebra, and thus they get a special name.

Definition 5.72 (Linear functional). 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 vector space $\mathcal{L}(V, \mathbf{F})$ also gets a special name and special notation.

Definition 5.73 (Dual space). The *dual space* of V is the vector space of linear functionals on V ; that is, $V^* := \mathcal{L}(V, \mathbf{F})$.

Lemma 5.74. Suppose V is finite-dimensional. Then V^* is also finite-dimensional, and

$$\dim V^* = \dim V.$$

Proof. By , we have

$$\dim V^* := \dim \mathcal{L}(V, \mathbf{F}) = (\dim V)(\dim \mathbf{F}) = \dim V$$

as desired. □

Definition 5.75 (Dual basis). If $\{v_1, \dots, v_n\}$ is 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

$$\phi_i(v_j) = \begin{cases} 1 & (i = j) \\ 0 & (i \neq j) \end{cases}$$

The following result states that dual basis gives coefficients for linear combination.

Proposition 5.76. Suppose $\{v_1, \dots, v_n\}$ is a basis of V , and $\{\phi_1, \dots, \phi_n\}$ is the dual basis. Then for each $v \in V$,

$$v = \phi_1(v)v_1 + \dots + \phi_n(v)v_n.$$

The following result states that the dual basis is a basis of the dual space.

Proposition 5.77. Suppose V is finite-dimensional. Then the dual basis of a basis of V is a basis of V^* .

Definition 5.78 (Dual map). Suppose $T \in \mathcal{L}(V, W)$. The **dual map** of T is the linear map $T^* \in \mathcal{L}(V, W)$ defined for each $\phi \in W^*$ by

$$T^*(\phi) = \phi \circ T.$$

Proposition 5.79 (Algebraic properties of dual map). Suppose $T \in \mathcal{L}(V, W)$. Then

- (1) $(S + T)^* = S^* + T^*$ for all $S \in \mathcal{L}(V, W)$
- (2) $(\lambda T)^* = \lambda T^*$ for all $\lambda \in \mathbf{F}$
- (3) $(ST)^* = T^*S^*$ for all $S \in \mathcal{L}(V, W)$

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 $\ker T$.

Definition 5.80 (Annihilator). For $U \subset V$, the **annihilator** of U is defined by

$$U^\circ := \{\phi \in V^* \mid \phi(u) = 0, \forall u \in U\}.$$

Proposition 5.81. $U^\circ \leq V$.

Proposition 5.82 (Dimension of annihilator). Suppose V is finite-dimensional, $U \leq V$. Then

$$\dim U^\circ = \dim V - \dim U.$$

The following are conditions for the annihilator to equal $\{0\}$ or the whole space.

Proposition 5.83. Suppose V is finite-dimensional, $U \leq V$. Then

- (i) $U^\circ = \{0\} \iff U = V$;
- (ii) $U^\circ = V^* \iff U = \{0\}$.

The following result concerns $\ker T^*$.

Proposition 5.84. Suppose V and W are finite-dimensional, $T \in \mathcal{L}(V, W)$. Then

- (i) $\ker T^* = (\operatorname{im} T)^\circ$;
- (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 5.85. *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 5.86. *Suppose V and W finite-dimensional, $T \in \mathcal{L}(V, W)$. Then*

- (i) $\dim \text{im } T^* = \dim \text{im } T$;
- (ii) $\dim T^* = (\ker T)^\circ$.

Proposition 5.87. *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 5.88. *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 ([Ax124] 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 ([Ax124] 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 ([Ax124] 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 ([Ax124] 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 ([Ax124] 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 ([Ax124] 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 ([Ax124] 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 ([Ax124] 3D). Suppose $T \in \mathcal{L}(V, W)$ is invertible. Show that T^{-1} is invertible and

$$(T^{-1})^{-1} = T.$$

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 $(T^{-1})^{-1} = T$. □

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). Suppose $U \leq V$, V/U is finite-dimensional. Prove that $V \cong U \times (V/U)$.

Solution.

$$\dim V = \dim U + (\dim V - \dim U) = \dim U + \dim(V/U).$$

□

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.

6 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 .

§6.1 Ordered Sets and Boundedness

Definitions

Let S be a set.

Definition 6.1. 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 6.2 (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 $x \in E$; β is called an *upper bound* of E .
- (ii) E is *bounded below* if there exists $\beta \in S$ such that $x \geq \beta$ for all $x \in E$; β is called a *lower bound* of E .

E is *bounded* in S if it is bounded above and below.

Definition 6.3 (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 6.4 (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 6.5. 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 6.6. 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 6.7. *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 6.8. *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.

Proposition 6.9 (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. □

For the rest of this section, suppose S has the least-upper-bound property.

Proposition 6.10 (Additive property). *Given non-empty subsets $A, B \subset S$, 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 Proposition 6.9 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 $a + b \leq c$.

□

Proposition 6.11 (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. □

Proposition 6.12. *Let $B \subset S$ be non-empty and bounded below. Let $A = \{-b \mid b \in B\}$. 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 . □

Ordered Fields

Definition 6.13 (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.

Proposition 6.14 (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 $\frac{1}{y} < \frac{1}{x}$.

□

§6.2 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 6.15. $x^2 - 2 = 0$ 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 irrational number $\sqrt{2}$.

Continuing from the above lemma,

Lemma 6.16. 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 6.17. \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 6.18 (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).*

Proof. We prove by construction, as follows. \square

We now want to construct \mathbb{R} from \mathbb{Q} ; one method to do so is using Dedekind cuts¹.

Definition 6.19 (Dedekind cut). A *Dedekind cut* $\alpha \subset \mathbb{Q}$ satisfies the following properties:

- (i) $\alpha \neq \emptyset, \alpha \neq \mathbb{Q}$;
- (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 6.20. 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 6.21. $\sqrt[3]{2}$ is not rational, but it is real. $\sqrt[3]{2}$ corresponds to the cut

$$\alpha = \{p \in \mathbb{Q} \mid p^3 < 2\}.$$

¹proposed by German mathematician Richard Dedekind in 1872.

- (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 6.22. 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 6.23. \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 6.24. 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 6.25 (Addition). Given $\alpha, \beta \in \mathbb{R}$,

$$\alpha + \beta := \{r \in \mathbb{Q} \mid r = a + b, a \in \alpha, b \in \beta\}.$$

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 $r < r'$.

□

Proposition 6.26.

- (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 $-\alpha$.

Proposition 6.27. 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 6.28. For all $\alpha, \beta \in \mathbb{R}^+$, we 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

$$\alpha\beta = \begin{cases} (-\alpha)(-\beta) & a < 0^*, \beta < 0^*, \\ -[(-\alpha)\beta] & a < 0^*, \beta > 0^*, \\ -[\alpha(-\beta)] & \alpha > 0^*, \beta < 0^*. \end{cases}$$

where we make negative numbers positive, multiply, and then negate them as needed.

Proposition 6.29.

- (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 6.30. *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} .

Properties of \mathbb{R}

Proposition 6.31 (\mathbb{R} is archimedean). *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 . □

Corollary 6.32. *Let $\varepsilon > 0$. Then there exists $n \in \mathbb{N}$ such that $0 < \frac{1}{n} < \varepsilon$.*

Proof. Take $x = \varepsilon$ and $y = 1$. □

Proposition 6.33 (\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; that is, construct the required p from the given x and y .

Since $x < y$, we have $y - x > 0$. By the archimedian property, there exists $n \in \mathbb{N}$ such that

$$\frac{1}{n} < y - x.$$

Consider the set comprising multiples of $\frac{1}{n}$. Since this set is unbounded, choose the first multiple $m \in \mathbb{N}$ such that $\frac{m}{n} > x$.

We now claim 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 6.34 (\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 $y^n = x$.*

Proof. The uniqueness of such y is clear, since $0 < y_1 < y_2$ implies $y_1^n < y_2^n$.

Claim. $y = \sup E$, where

$$E = \{t \in \mathbb{R}^+ \mid t^n < x\}.$$

We first show that E has a supremum, by showing that it is (i) non-empty, and (ii) bounded above:

- (i) Let $t = \frac{x}{1+x}$. Then $0 \leq t < 1$, so $t^n \leq t < x$. Hence $t \in E$, which implies $E \neq \emptyset$.
- (ii) We claim that $1+x$ is an upper bound for E . If $t > 1+x$ then $t^n \geq t > x$, so that $t \notin E$. Hence $1+x$ is an upper bound of E .

Hence E has a supremum; let $y = \sup E$.

To prove that $y^n = x$, we show that the inequalities (i) $y^n < x$ and (ii) $y^n > x$ lead to a contradiction. Consider the identity $b^n - a^n = (b-a)(b^{n-1} + b^{n-2}a + \cdots + a^{n-1})$, which yields the inequality

$$b^n - a^n < (b-a)nb^{n-1}$$

when $0 < a < b$.

- (i) Suppose $y^n < x$. Choose h so that $0 < h < 1$ and

$$h < \frac{x - y^n}{n(y+1)^{n-1}}.$$

let $a = y, b = y + h$. Then

$$(y+h)^n - y^n < hn(y+h)^{n-1} < hn(y+1)^{n-1} < x - y^n.$$

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 .

(ii) Suppose $y^n > x$. Let

$$k = \frac{y^n - x}{ny^{n-1}}.$$

Then $0 < k < y$. If $t \geq y - k$, we conclude that

$$y^n - t^n \leq y^n - (y - k)^n < kny^{n-1} = y^n - x.$$

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 .

Hence $y^n = x$, and the proof is complete. □

Notation. y is denoted by $\sqrt[n]{x}$ or $x^{\frac{1}{n}}$.

Corollary 6.35. *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

$$(ab)^{\frac{1}{n}} = \alpha\beta = a^{\frac{1}{n}} b^{\frac{1}{n}}.$$

□

The next result shows that real numbers can be approximated to any desired degree of accuracy by rational numbers with finite decimal representations.

Proposition 6.36. *Assume $x \geq 0$. 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; that is, we 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

$$a_0 + \frac{a_1}{10} + \dots + \frac{a_n}{10^n} \leq a_0 + \frac{a_1}{10} + \dots + \frac{a_n + 1}{10^n}.$$

Then $0 \leq a_n \leq 9$ and we have

$$r_n \leq x < r_n + \frac{1}{10^n},$$

where $r_n = a_0.a_1a_2 \dots a_n$. □

Furthermore, it is easy to verify that $x = \sup_{n \in \mathbb{N}} r_n$.

Extended Real Number System

Definition 6.37 (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 6.38. *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*.

§6.3 Complex Field

Consider the Cartesian product \mathbb{R}^2 . A *complex number* is an ordered pair $(a, b) \in \mathbb{R}^2$.

Proposition 6.39. *Let $x = (a, b)$, $y = (c, d)$ be two complex numbers. We write $x = y$ if and only if $a = c$ and $b = d$. \mathbb{R}^2 , with addition and multiplication defined as*

$$\begin{aligned}x + y &= (a + c, b + d) \\ xy &= (ac - bd, ad + bc)\end{aligned}$$

*is a field. Note that the additive identity is $(0, 0)$, and multiplicative identity is $(1, 0)$. We call this structure \mathbb{C} , the **complex field**.*

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} .

Proposition 6.40. *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 $i = (0, 1)$.

Proposition 6.41. $i^2 = -1$.

Proof.

$$i^2 = (0, 1)(0, 1) = (-1, 0) = -1.$$

□

Proposition 6.42. *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 $a, b \in \mathbb{R}$, $z = a + bi$, 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 .

Proposition 6.43. For $z, w \in \mathbb{C}$,

- (i) $\overline{z + w} = \bar{z} + \bar{w}$
- (ii) $\overline{z\bar{w}} = \bar{z} w$
- (iii) $z + \bar{z} = 2 \operatorname{Re}(z)$, $z - \bar{z} = 2i \operatorname{Im}(z)$
- (iv) $z\bar{z} \in \mathbb{R}$ 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 6.44. For $z, w \in \mathbb{C}$,

- (i) $|z| \geq 0$
- (ii) $|\bar{z}| = |z|$
- (iii) $|zw| = |z||w|$
- (iv) $|\operatorname{Re}(z)| \leq |z|$

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

$$|\operatorname{Re}(z)| = |a| = \sqrt{a^2} \leq \sqrt{a^2 + b^2} = |z|.$$

□

Theorem 6.45 (Triangle inequality). For $z, w \in \mathbb{C}$,

$$|z + w| \leq |z| + |w|. \quad (6.1)$$

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 6.46 (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 6.47. 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 $y - z$.

□

Theorem 6.48 (Cauchy–Schwarz inequality). *If $a_1, \dots, a_n, b_1, \dots, b_n \in \mathbb{C}$, then*

$$\left| \sum_{i=1}^n a_i \bar{b}_i \right|^2 \leq \sum_{i=1}^n |a_i|^2 \sum_{i=1}^n |b_i|^2. \quad (6.2)$$

Proof. Let

$$A = \sum_{i=1}^n |a_i|^2, \quad B = \sum_{i=1}^n |b_i|^2, \quad C = \sum_{i=1}^n a_i \bar{b}_i.$$

If $B = 0$, then $b_1 = \dots = b_n = 0$, and the conclusion is trivial. Now assume that $B > 0$. Then consider the sum

$$\begin{aligned} \sum_{i=1}^n |Ba_i - Cb_i|^2 &= \sum_{i=1}^n (Ba_i - Cb_i)(\overline{Ba_i - Cb_i}) \\ &= \sum_{i=1}^n (Ba_i - Cb_i)(B\bar{a}_i - \overline{Cb_i}) \\ &= B^2 \sum_{i=1}^n |a_i|^2 - B\bar{C} \sum_{i=1}^n a_i \bar{b}_i - BC \sum_{i=1}^n \bar{a}_i b_i + |C|^2 \sum_{i=1}^n |b_i|^2 \\ &= B^2 A - B|C|^2 \\ &= B(AB - |C|^2). \end{aligned}$$

Since each term $\sum_{i=1}^n |Ba_i - Cb_i|^2$ is non-negative, we have that $\sum_{i=1}^n |Ba_i - Cb_i|^2 \geq 0$, and so

$$B(AB - |C|^2) \geq 0.$$

Since $B > 0$, it follows that $AB - |C|^2 \geq 0$. 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 $\mathbf{a} \in \mathbb{C}^n$:

$$|\mathbf{a}| = \langle \mathbf{a}, \mathbf{a} \rangle^{\frac{1}{2}}.$$

§6.4 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 6.49. 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 6.50 (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) $\|\mathbf{x} \cdot \mathbf{y}\| \leq \|\mathbf{x}\| \|\mathbf{y}\|$
- (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) This is an immediate consequence of the Cauchy–Schwarz inequality.

(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 6.1 ([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 6.2 ([Rud76] 1.2). Prove that there is no rational number whose square is 12.

Solution. Prove by contradiction. \square

Exercise 6.3 ([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 6.4 ([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 Proposition 6.14, an order $<$ that makes \mathbb{C} an ordered field would have to satisfy $-1 = i^2 > 0$, contradicting $1 > 0$. \square

Exercise 6.5 ([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 6.6 ([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

$$2ab = (|w| + u)^{\frac{1}{2}} (|w| - u)^{\frac{1}{2}} = (|w|^2 - u^2)^{\frac{1}{2}} = (v^2)^{\frac{1}{2}} = |v|.$$

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 6.7 ([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

7 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	$B_r(p) \setminus \{p\}$
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 7.1: Notation for Chapter 7

§7.1 Metric Spaces

Definitions and Examples

Definition 7.1 (Metric space). A **metric space** is a set X with an associated *metric* $d : X \times X \rightarrow \mathbb{R}$, which satisfies the following properties 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)

For the rest of the chapter, X is taken to be a metric space, unless specified otherwise.

Example 7.2 (Metrics on \mathbb{R}^n). Each of the following functions define metrics on \mathbb{R}^n .

$$\begin{aligned} d_1(x, y) &= \sum_{i=1}^n |x_i - y_i|; \\ d_2(x, y) &= \sqrt{\sum_{i=1}^n (x_i - y_i)^2}; \\ d_\infty(x, y) &= \max_{i \in \{1, 2, \dots, n\}} |x_i - y_i|. \end{aligned}$$

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 7.3 (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 7.4 (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:

$$B_r(p) \setminus \{p\} = \{q \in X \mid 0 < d(p, q) < r\}.$$

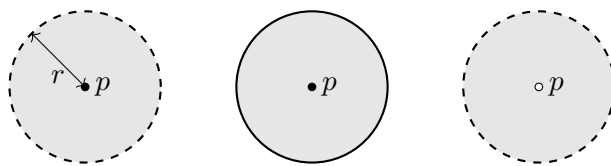


Figure 7.1: Open ball, closed ball, punctured ball

Example 7.5. 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 7.6 (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 $E \subset B_M(p)$.

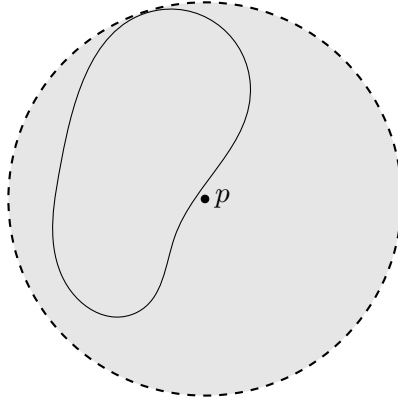


Figure 7.2: Bounded set

Proposition 7.7. 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

Definition 7.8 (Diameter). Let non-empty $E \subset X$. Then the **diameter** of E is

$$\text{diam } E := \sup_{p, q \in E} d(p, q).$$

Example 7.9. Find the diameter of the open unit ball in \mathbb{R}^n :

$$B = \{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\| < 1\}.$$

Proof. Note that for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,

$$d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| \leq \|\mathbf{x}\| + \|-\mathbf{y}\| = \|\mathbf{x}\| + \|\mathbf{y}\| < 1 + 1 = 2.$$

On the other hand, for any $\varepsilon > 0$, we pick

$$\mathbf{x} = \left(1 - \frac{\varepsilon}{4}, 0, \dots, 0\right), \quad \mathbf{y} = \left(-\left(1 - \frac{\varepsilon}{4}\right), 0, \dots, 0\right).$$

Then $d(\mathbf{x}, \mathbf{y}) = 2 - \frac{\varepsilon}{2} > 2 - \varepsilon$. Since ε is arbitrary, we have that $\text{diam } B = 2$. \square

Proposition 7.10. $E \subset \mathbb{R}^n$ is bounded if and only if $\text{diam } E < +\infty$.

Proof.

\Rightarrow If E is bounded, then there exists $M > 0$ such that $\|x\| \leq M$ for all $x \in E$.

Thus for any $x, y \in E$,

$$d(x, y) = \|x - y\| \leq \|x\| + \|y\| \leq 2M.$$

Thus $\text{diam } E = \sup d(x, y) \leq 2M < +\infty$.

\Leftarrow Suppose that $\text{diam } E = r$. Pick a random point $x \in E$, suppose that $\|x\| = R$.

Then for any other $y \in E$,

$$\|y\| = \|x + (y - x)\| \leq \|x\| + \|y - x\| \leq R + r.$$

Thus, by picking $M = R + r$, we obtain $\|y\| \leq M$ for all $y \in E$, and we are done. \square

Remark. Basically we used x to confine E within a ball, which is then confined within an even bigger ball centered at the origin.

Open and Closed Sets

Definition 7.11 (Neighbourhood). $N \subset X$ is a *neighbourhood* of $p \in X$ if there exists $\varepsilon > 0$ such that $B_\varepsilon(p) \subset N$.

Definition 7.12 (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 $B_\varepsilon(p) \subset E$.

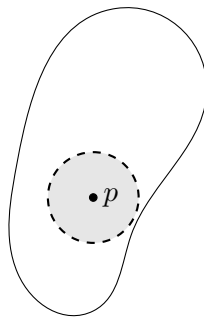


Figure 7.3: Open set

Proposition 7.13. *Any open ball is open.*

Proof. Let $B_r(p)$ be an open ball. Let $q \in B_r(p)$, then $d(p, q) < r$. Let $\varepsilon = r - d(p, q)$, and note that $\varepsilon > 0$.

Consider the ball $B_\varepsilon(q)$, and let $s \in B_\varepsilon(q)$. By the triangle inequality,

$$\begin{aligned} d(p, s) &\leq d(q, s) + d(p, q) \\ &< \varepsilon + d(p, q) \\ &= r \end{aligned}$$

and thus $s \in B_r(p)$. Since for all $q \in B_r(p)$ there exists $\varepsilon > 0$ such that $B_\varepsilon(q) \subset B_r(p)$, we have that $B_r(p)$ is open. \square

Proposition 7.14.

- (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 $B_\delta(p) \subset \bigcap_{i \in I} E_i$. \square

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 7.15. *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.

\Rightarrow We prove by construction; that is, construct the required set G .

¹notice that the definition of an open ball depends on the metric space!

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 Proposition 7.14, 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 $E = G \cap Y$.

$\boxed{\Leftarrow}$ 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 . \square

The complement of an open set is a closed set.

Definition 7.16 (Closed set). $E \subset X$ is **closed** if its complement $E^c = X \setminus E$ is open.

Proposition 7.17. Any closed ball is closed.

Proof. To prove that $\overline{B}_r(p)$ is closed, we need to show that its complement

$$\overline{B}_r(p)^c = \{q \in X \mid d(p, q) > r\}$$

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

Proposition 7.18.

- (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 Proposition 7.14, simply take complements and apply de Morgan's laws. \square

Remark. As above, the indexing set in (iii) must be finite; for instance, the closed intervals $F_i = \left[-1 + \frac{1}{i}, 1 - \frac{1}{i}\right]$ are all closed in \mathbb{R} , but their union $\bigcup_{i=1}^{\infty} F_i = (-1, 1)$ is open.

Interior, Closure, Boundary

Definition 7.19. Suppose $E \subset X$.

- (i) The **interior** of E , denoted by E° , is the union of all open subsets of X contained in E ; $p \in E^\circ$ is an **interior point** of E . (Equivalently, E° is the set of all points in E for which E is a neighbourhood; p is an interior point if there exists $\varepsilon > 0$ such that $B_\varepsilon(p) \subset E$.)
- (ii) The **closure** of E , denoted by \overline{E} , is the intersection of all closed subsets of X containing E . E is said to be **dense** if $\overline{E} = X$. (Equivalently, every point of X is either a limit point of E , or in E .)
- (iii) The **boundary** of E is $\partial E = \overline{E} \setminus E^\circ$; $p \in \partial E$ is 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.

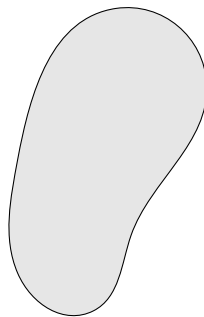


Figure 7.4: Interior, closure, boundary

Example 7.20. The interior of the closed interval $[a, b] \subset \mathbb{R}$ is the open interval (a, b) . \mathbb{Q} is dense in \mathbb{R} .

Remark. It is obvious that the interior E° is open since the union of open sets is open; similarly, the closure \overline{E} is closed since the intersection of closed sets is closed.

Proposition 7.21. Suppose $E \subset X$.

- (i) E is open if and only if $E = E^\circ$.
- (ii) E is closed if and only if $E = \overline{E}$.

Proof.

- (i) \Rightarrow Suppose E is open. Then E is an open subset of X contained in E (since $E \subset E$), so $E \subset E^\circ$. Conversely, suppose $A \subset E^\circ$. Then A is an open subset of X contained in E , so $A \subset E$, and hence $E^\circ \subset E$. Therefore $E = E^\circ$.

\Leftarrow Since an arbitrary union of open sets is open, E° is itself an open set. Clearly E° is the unique largest open subset of X contained in E , since all open subsets of E are contained in E° .

- (ii) \Rightarrow If E is itself closed then evidently $E = \overline{E}$.

\Leftarrow Since an arbitrary intersection of closed sets is closed, \overline{E} is the unique smallest closed subset of X containing E .

□

Proposition 7.22. 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.

\Rightarrow Suppose that $p \in \overline{E}$. Suppose, for a contradiction, that there exists some 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.

\Leftarrow Suppose that every ball $B_\varepsilon(p)$ meets E . Suppose, for a contradiction, that $p \notin \overline{E}$. Then 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 .

Proposition 7.23. Suppose $E \subset X$. Let $F \supset E$ be some closed set. Then $\overline{E} \subset F$.

Proof. Let p be a limit point of E . Then p is a limit point of F . But since F is closed, F contains all its limit points, so all the limit points of E are in F . Hence $\overline{E} \subset F$. □

Remark. This means that \overline{E} is the “smallest” closed set containing E .

Limit Points

Definition 7.24.

- (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** (or *accumulation 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 $B_\varepsilon(p) \cap E = \{p\}$).

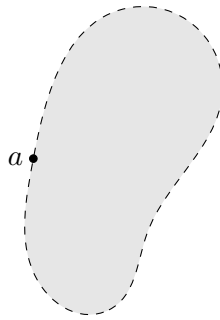


Figure 7.5: Adherent point, limit point, isolated point

Example 7.25 (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 7.26 (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 the closed interval $[a, b]$ is a limit point of the set of numbers in the open interval (a, b) .
- Consider the metric space \mathbb{R}^2 . The limit point set 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 7.27. *If p is a limit point of E , then every ball of p contains infinitely many points of E .*

Proof. We prove by contradiction. Suppose otherwise, for a contradiction, that there exists $B_r(p)$ which contains only a finite number of points of E distinct from p . Then 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 7.28. *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 (Theorem 7.50).

A closed set was defined to be the complement of an open set. The next result describes closed sets in another way.

Proposition 7.29. *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 prove that $p \in E$.

Suppose otherwise, 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

Proposition 7.30. *Suppose $E \subset X$. Then E' is a closed subset of X .*

Proof. To prove that E' is closed, we need to show that its complement $(E')^c$ is open.

Suppose $p \in (E')^c$. Then p is not a limit point of E , so there exists a ball $B_\varepsilon(p)$ whose intersection with E is either empty or $\{p\}$ (depending on whether $p \in E$ or not).

Claim. $B_{\frac{\varepsilon}{2}}(p) \subset (E')^c$.

Let $q \in B_{\frac{\varepsilon}{2}}(p)$.

- If $q = p$, then clearly $q \in (E')^c$.
- If $q \neq p$, 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

Proposition 7.31. *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 .

$\boxed{\overline{E} \subset E \cup E'}$ 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. \square

Proposition 7.32. *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 $y \in \overline{E}$. \square

(y is either in E , or a limit point of E)

review
proof

§7.2 Compactness

Definitions and Properties

Definition 7.33 (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 7.34 (Compactness). $K \subset X$ is **compact** if *every* open cover of K has 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

$$K \subset \bigcup_{k=1}^n U_{i_k}.$$

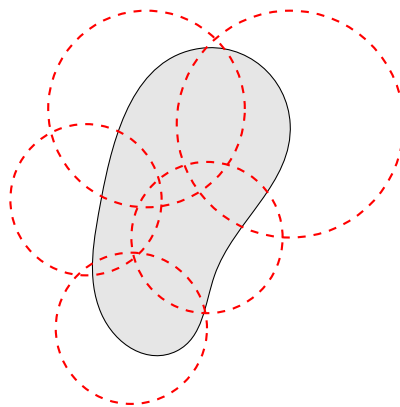


Figure 7.6: Compact set

Example 7.35. • The real line \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 Proposition 7.40 for the proof.)

Proposition 7.36. *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 need to 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} . \square

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 7.37. *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.

\Rightarrow Suppose K is compact relative to X . We will show that K is compact relative to Y . 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 for all $i \in I$, U_i is open relative to Y , by Proposition 7.15, there exists V_i open relative to X such that $U_i = Y \cap V_i$. Consider $\{V_i \mid i \in I\}$, which is an open cover of K , since it is a collection of open sets. 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 .

\Leftarrow 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 $U_i \subset V_i$. \square

Proposition 7.38. *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. \square

Proposition 7.39. *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 taking the union,

$$\begin{aligned} \bigcup_{k=1}^n (B_\varepsilon(p) \cap B_{r_{i_k}}(q_{i_k})) &= \emptyset \\ B_\varepsilon(p) \cap \left(\bigcup_{k=1}^n B_{r_{i_k}}(q_{i_k}) \right) &= \emptyset \\ B_\varepsilon(p) \cap W &= \emptyset \end{aligned}$$

as desired. □

Proposition 7.40. *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 want to show that there exists 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. \square

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 Theorem 7.49).

Corollary 7.41. *If F is closed and K is compact, then $F \cap K$ is compact.*

Proof. Suppose F is closed, K is compact. By Proposition 7.39, K is closed. By Proposition 7.18, 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 Proposition 7.40, $F \cap K$ is compact. \square

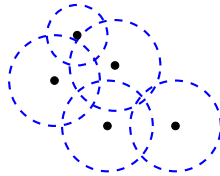
Heine–Borel Theorem

Proposition 7.42. *K is compact if and only if every infinite subset of K has a limit point in K .*

Proof.

\Rightarrow Suppose K is compact. Let E be an infinite subset of K . Suppose otherwise, 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 .

\Leftarrow Suppose every infinite subset of K 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 $p \in U_p$.

\square

To complete
proof

Proposition 7.43 (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. For $n \in \mathbb{N}$, let $I_n = [a_n, b_n]$. 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 7.44. 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 7.45. Suppose (I_n) is a decreasing sequence of k -cells; that is, $I_1 \supset I_2 \supset \dots$. Then $\bigcap_{n=1}^{\infty} I_n \neq \emptyset$.

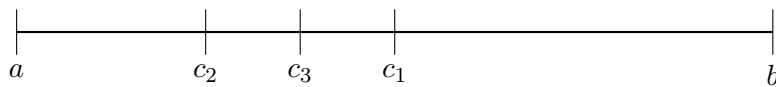
Proof. For $n \in \mathbb{N}$, let

$$I_n = \{(x_1, \dots, x_k) \mid a_{ni} \leq x_i \leq b_{ni}, 1 \leq i \leq k\},$$

and let $I_{n_i} = [a_{n_i}, b_{n_i}]$, where $I_n = I_{n_1} \times \dots \times I_{n_k}$ is the Cartesian product of k closed intervals.

For each i ($i = 1, \dots, k$), □

Lemma 7.46. Every closed interval is compact (in \mathbb{R}).



Proof. Suppose otherwise, 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

$$I_1 \supset I_2 \supset I_3 \supset \cdots$$

where all of them have no finite subcover of \mathcal{U} .

By the nested interval theorem (Proposition 7.43), 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 7.47. *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

$$\delta = \left(\sum_{i=1}^k (b_i - a_i)^2 \right)^{\frac{1}{2}}$$

that is, δ 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 . Then

$$|\mathbf{x} - \mathbf{y}| \leq \delta \quad (\forall \mathbf{x}, \mathbf{y} \in I)$$

Suppose otherwise, 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) $I \supset I_1 \supset I_2 \supset \cdots$
- (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 Theorem 2.39, 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 7.48. *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 . From Lemma 7.47, we have that I is compact. Since E is a closed subset of a compact set, by Proposition 7.40, E is compact.

(ii) \implies (iii) This directly follows from Proposition 7.42.

(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. \square

review
proof

Theorem 7.49 (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. \square

Bolzano–Weierstrass Theorem

Theorem 7.50 (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 Proposition 7.42, E has a limit point in I and thus \mathbb{R}^n . \square

Cantor Intersection Theorem

A collection $\mathcal{A} = \{A_i \mid i \in I\}$ 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 7.51. *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 otherwise, 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} . \square

Theorem 7.52 (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 follows from the previous result; it is obvious that the intersection of every finite subcollection of a decreasing sequence of sets must be non-empty. \square

The following result is a characterisation of compact sets.

Proposition 7.53. *K is compact if and only if every collection of closed subsets of K satisfies the finite intersection property.*

Proof.

\Rightarrow 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 7.54 (Sequential compactness). $K \subset X$ 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 7.55. *$K \subset X$ 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 $\{U_\alpha \mid \alpha \in \Lambda_\alpha\}$ 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

§7.3 Perfect Sets

Definition 7.56 (Perfect set). E is *perfect* if

- (i) E is closed;
- (ii) every point of E is a limit point of E .

Proposition 7.57. *Let non-empty $P \subset \mathbb{R}^k$ be perfect. Then P is uncountable.*

Proof. Since P has limit points, P must be infinite. Suppose, for a contradiction, that P is countable. This means we can list the points of P in a sequence:

$$\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots$$

We now construct 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 a neighborhood V_{n+1} such that (i) $V_{n+1} \subset B_n$, (ii) $\mathbf{x}_n \notin V_{n+1}$, (iii) $V_{n+1} \cap P$ is not empty. By (iii), V_{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 $\overline{B}_n \setminus K_{n+1}$. Since $K_n \subset P$, this implies that $\overline{B}_n \setminus K_{n+1}$ is empty. But each K_n is nonempty, by (iii), and $K_n \supset K_{n+1}$, by (i); this contradicts the Corollary to Theorem 2.36. \square

Corollary 7.58. *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

$$C_2 = \left[0, \frac{1}{9}\right] \cup \left[\frac{2}{9}, \frac{3}{9}\right] \cup \left[\frac{6}{9}, \frac{7}{9}\right] \cup \left[\frac{8}{9}, 1\right].$$

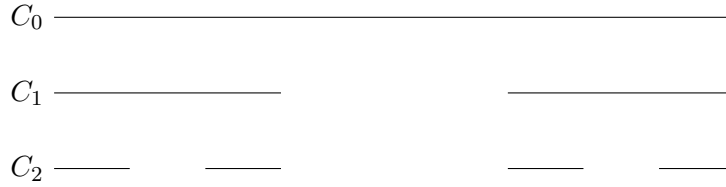


Figure 7.7: Cantor set

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 $C_{n+1} = \frac{1}{3}C_n \cup \left(\frac{1}{3}C_n + \frac{2}{3}\right)$.

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 7.59 (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 7.60. *C is a perfect set in \mathbb{R} which contains no segment.*

Proof. We prove (i) C contains no segment, and (ii) C is perfect.

(i) No segment 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 segment (α, β) contains a segment of the above form, if

$$3^{-m} < \frac{\beta - \alpha}{6},$$

C contains no segment.

(ii) Since we have shown that C is closed, it suffices to show that every point of C is a limit point of C .

Let $x \in C$, and let S be any segment 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 7.61. *C is uncountable.*

The following are a few more interesting properties of the Cantor set.

Proposition 7.62. *C is precisely the set of all real numbers in $[0, 1]$ whose ternary expansion contain only 0's or and 2's.*

Remark. Finite decimal expansions, as always, are not formally well-defined; for instance, $\frac{1}{3} = 0.1 = 0.0222\ldots$ so that $\frac{1}{3} \in C$ because it can be expressed with only zeros and twos in at least one of its ternary expansions.

Proof. □

Proposition 7.63. *C has measure zero; that is, for all $\varepsilon > 0$, C can be covered by intervals of total length less than ε .*

§7.4 Connectedness

Definition 7.64 (Connectedness). A and B are *separated* if

- (i) $A \cap \overline{B} = \emptyset$, and
- (ii) $\overline{A} \cap B = \emptyset$;

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, the interval $[0, 1]$ and the segment $(1, 2)$ are not separated, since 1 is a limit point of $(1, 2)$. However, the segments $(0, 1)$ and $(1, 2)$ are separated.

Example 7.65. 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 $A \cap \overline{B} = \emptyset$ and $\overline{A} \cap B = \emptyset$.

Proposition 7.66. Closed intervals in \mathbb{R} are connected.

Proof. Suppose otherwise, for a contradiction, that a closed interval $[a, b]$ is not connected. Then by definition, there exists non-empty sets A and B , with $A \cap \overline{B} = \emptyset$ and $\overline{A} \cap B = \emptyset$. WLOG let $a \in A$.

Let $s = \sup A$. Then by Proposition 7.32, $s \in \overline{A}$. Then $\overline{A} \cap B = \emptyset$ implies $s \notin B$, so $s \in A$. Thus $A \cap \overline{B} = \emptyset$ implies $s \notin \overline{B}$. 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

Proposition 7.67. $E \subset \mathbb{R}$ 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 Theorem 2.28, $z \in \overline{A}$; hence $z \notin B$. In particular, $x \leq z < y$.

If $z \notin A$, it follows that $x < z < y$ and $z \notin E$.

If $z \in A$, 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 $z_1 \notin E$. \square

Proposition 7.68. *The Cantor set C is totally disconnected.*

Exercises

Exercise 7.1. 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 7.2 (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

8 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.

§8.1 Sequences

Convergence

A **sequence** (x_n) in X is a function $f : \mathbb{N} \rightarrow X$ which maps $n \mapsto x_n$.

The *range* of a sequence (x_n) is the set

$$\{a \in X \mid \exists n \in \mathbb{N}, a = x_n\}.$$

Note that the range of a sequence may be a finite set or it may be infinite. (x_n) is *bounded* if its range is bounded.

Definition 8.1. A sequence (x_n) **converges** to $x \in X$, denoted by $x_n \rightarrow x$, if

$$\forall \varepsilon > 0, \quad \exists N \in \mathbb{N}, \quad \forall n \geq N, \quad d(x_n, x) < \varepsilon.$$

We call x a *limit* of (x_n) . If (x_n) does not converge, it is said to *diverge*.

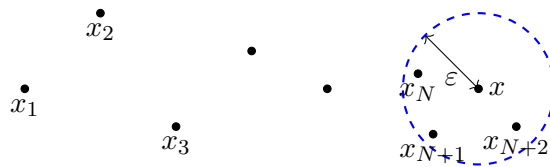


Figure 8.1: Convergence of sequence

Remark. This limit process conveys the intuitive idea that x_n can be made arbitrarily close to x , provided that n is sufficiently large. (Equivalently, if we remove more and more initial terms from the sequence, we see that the *tail* of the sequence should be increasingly closer to x .)

Remark. If $x_n \not\rightarrow x$, simply negate the definition for convergence:

$$\exists \varepsilon > 0, \quad \forall N \in \mathbb{N}, \quad \exists n \geq N, \quad d(x_n, x) \geq \varepsilon.$$

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 8.2. $\frac{1}{n} \rightarrow 0$.

Proof. Fix $\varepsilon > 0$. By the Archimedian 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 8.3. 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

$$\begin{aligned} |a_n - 1| &< \varepsilon \\ \frac{1}{\sqrt{n}} &= \left| (-1)^n \frac{1}{\sqrt{n}} \right| < \varepsilon \\ \frac{1}{n} &< \varepsilon^2 \\ n &> \frac{1}{\varepsilon^2} \end{aligned}$$

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 8.4 (Uniqueness of limit). *If a sequence converges, then its limit is unique.*

Proof. Let (x_n) be a sequence in X . Suppose that $x_n \rightarrow x$ and $x_n \rightarrow x'$ for $x, x' \in X$. We will show that $x' = x$.

Let $\varepsilon > 0$ be given. Then there exists $N, N' \in \mathbb{N}$ such that

$$\begin{aligned} n \geq N &\implies d(x_n, x) < \frac{\varepsilon}{2} \\ n \geq N' &\implies d(x_n, x') < \frac{\varepsilon}{2} \end{aligned}$$

Take $N_1 := \max\{N, N'\}$. If $n \geq N_1$, then both hold. By the triangle inequality,

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

Since this holds for all $\varepsilon > 0$, we must have $d(x, x') = 0$. Hence $x = x'$. \square

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

Notation. If (x_n) converges to x , we denote $\lim_{n \rightarrow \infty} x_n = x$.

We now outline some important properties of convergent sequences in metric spaces.

Proposition 8.5. *Let (x_n) be a sequence in X .*

- (i) $x_n \rightarrow x$ if and only if every open ball of x contains x_n for all but finitely many n .
- (ii) If (x_n) converges, then (x_n) is bounded.
- (iii) Suppose $E \subset X$. Then x is a limit point of E if and only if there exists a sequence (x_n) in $E \setminus \{x\}$ such that $x_n \rightarrow x$.

Proof.

- (i) \implies Suppose $x_n \rightarrow x$. Let $\varepsilon > 0$ be given, then there exists $N \in \mathbb{N}$ such that

$$n \geq N \implies d(x_n, x) < \varepsilon.$$

Corresponding to this ε , consider the open ball $B_\varepsilon(x)$. Then by definition, for $y \in X$,

$$d(y, x) < \varepsilon \implies y \in B_\varepsilon(x).$$

Hence $n \geq N$ implies $x_n \in B_\varepsilon(x)$.

\impliedby Suppose every open ball of x contains all but finitely many of the x_n .

Let $\varepsilon > 0$ be given. Consider the open ball $B_\varepsilon(x)$. Since $B_\varepsilon(x)$ is a open ball of x , it will also eventually contain all x_n ; that is, there exists $N \in \mathbb{N}$ such that if $n \geq N$, then $x_n \in B_\varepsilon(x)$, i.e. $d(x_n, x) < \varepsilon$. Hence $x_n \rightarrow x$.

- (ii) Suppose $x_n \rightarrow x$. Let $\varepsilon > 0$ be given. Then there exists $N \in \mathbb{N}$ such that $n \geq N$ implies $d(x_n, x) < 1$. Now let

$$r = \max\{1, d(x_1, x), \dots, d(x_N, x)\}.$$

Then $d(x_n, x) \leq r$ for $n = 1, 2, \dots, N$, so the range of x_n is bounded by $B_r(x)$. Hence (x_n) is bounded.

- (iii) \implies Suppose x is a limit point of E .

Consider a sequence of open balls $\left(B_{\frac{1}{n}}(x)\right)$, for $n \in \mathbb{N}$. Since x is a limit point, each open ball intersects with E at some point which is not x . We pick one such point x_n from each $B_{\frac{1}{n}}(x) \cap E$. Then

$$d(x_n, x) < \frac{1}{n}.$$

Let $\varepsilon > 0$ be given. Then by the Archimedian property, there exists $N \in \mathbb{N}$ such that $\frac{1}{N} < \varepsilon$. If $n \geq N$,

$$d(x_n, x) \leq \frac{1}{n} \leq \frac{1}{N} < \varepsilon,$$

which shows that $x_n \rightarrow x$.

$\boxed{\Leftarrow}$ Suppose that there exists a sequence (x_n) in $E \setminus \{x\}$ such that $x_n \rightarrow x$. Then for each open ball $B_\varepsilon(x)$, we can find some $N \in \mathbb{N}$ such that if $n \in \mathbb{N}$ then

$$x_n \in B_\varepsilon(x).$$

Since $x_n \in E \setminus \{x\}$, this shows that x is a limit point of E .

□

Proposition 8.6 (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.

Proposition 8.7 (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*

$$(i) \quad \lim_{n \rightarrow \infty} ca_n = ca, \text{ where } c \text{ is a constant} \quad (\text{scalar multiplication})$$

$$(ii) \quad \lim_{n \rightarrow \infty} (a_n + b_n) = a + b \quad (\text{addition})$$

$$(iii) \lim_{n \rightarrow \infty} (a_n b_n) = ab \quad (\text{multiplication})$$

$$(iv) \lim_{n \rightarrow \infty} \frac{a_n}{b_n} = \frac{a}{b} \quad (b_n \neq 0, b \neq 0) \quad (\text{division})$$

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 $n \geq N$,

$$\left| \frac{1}{b_n} - \frac{1}{b} \right| = \left| \frac{b - b_n}{b_n b} \right| < \frac{2}{|b|^2} |b_n - b| < \varepsilon.$$

□

Proposition 8.8 (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 $c_n \rightarrow L$. □

Given a complex sequence (z_n) , there are two associated real sequences $(\operatorname{Re}(z_n))$ and $(\operatorname{Im}(z_n))$. The next result relates convergence of (z_n) to convergence of $(\operatorname{Re}(z_n))$ and $(\operatorname{Im}(z_n))$.

Proposition 8.9 (Complex sequences). *Let (z_n) be a complex sequence. Write $z_n = x_n + iy_n$ with $x_n, y_n \in \mathbb{R}$, so that (x_n) and (y_n) are real sequences. Then (z_n) converges if and only if both (x_n) and (y_n) converge. Moreover, in the case where (z_n) converges, we have*

$$\lim_{n \rightarrow \infty} z_n = \lim_{n \rightarrow \infty} x_n + i \lim_{n \rightarrow \infty} y_n.$$

Subsequences

Definition 8.10 (Subsequence). Given a sequence (x_n) , consider a sequence (n_k) of positive integers such that $n_1 < n_2 < \dots$. Then (x_{n_k}) is called a **subsequence** of (x_n) . If (x_{n_k}) converges, its limit is called a *subsequential limit* of (x_n) .

Proposition 8.11. (x_n) converges to x if and only if every subsequence of (x_n) converges to x .

Proof.

\Rightarrow Suppose $x_n \rightarrow x$. Let $\varepsilon > 0$ be given. Then there exists $N \in \mathbb{N}$ such that

$$n \geq N \implies d(x_n, x) < \varepsilon.$$

Every subsequence of (x_n) can be written in the form (x_{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

$$k > M \implies n_k > n_M \geq N \implies d(x_{n_k}, x) < \varepsilon.$$

Hence every subsequence of (x_n) converges to x .

\Leftarrow Suppose every subsequence of (x_n) converges to x . Since (x_n) is a subsequence of itself, we must have $x_n \rightarrow x$. \square

Proposition 8.12. In a compact metric space, any sequence has a convergent subsequence.

Proof. Suppose (x_n) is a sequence in a compact metric space X .

Let E be the range of (x_n) . We have to consider two cases: (i) E is finite, (ii) E is infinite. For both cases, we will construct the desired convergent subsequence.

- (i) Notice that there are infinitely many terms in the sequence (x_n) , but only finitely many distinct terms in E . Hence by the pigeonhole principle, at least one term of E appears infinitely many times in the sequence.

That is, there exists $x \in E$ and a sequence (n_k) with $n_1 < n_2 < \dots$ such that

$$x_{n_1} = x_{n_2} = \dots = x.$$

This subsequence (x_{n_k}) that we have constructed evidently converges to x .

- (ii) If E is infinite, then E is an infinite subset of a compact set. By Proposition 7.42, E has a limit point $x \in X$.

We now construct a subsequence (x_{n_k}) of (x_n) such that $x_{n_k} \rightarrow x$. Choose n_1 so that $d(x, x_{n_1}) < 1$. Having chosen n_1, \dots, n_{k-1} , choose n_k where $n_k > n_{k-1}$ such that $d(x, x_{n_k}) < \frac{1}{k}$ (such n_k exists due to Proposition 7.27). Then $x_{n_k} \rightarrow x$. \square

Corollary 8.13 (Bolzano–Weierstrass). *Every bounded sequence in \mathbb{R}^k has a convergent subsequence.*

Proof. By Proposition 7.48, 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 . \square

Lemma 8.14. *Suppose (x_n) is a sequence in X . Then the subsequential limits of (x_n) form a closed subset of X .*

Proof. Let E be the set of all subsequential limits of (x_n) , let q be a limit point of E . We want to show that $q \in E$.

Choose n_1 so that $x_{n_1} \neq q$. (If no such n_1 exists, then E has only one point, and there is nothing to prove.) Put $\delta = d(q, x_{n_1})$. Suppose n_1, \dots, n_{i-1} are chosen. Since q is a limit point of E , there is an $x \in E$ with $d(x, q) < 2^{-1}\delta$. Since $x \in E$, there is an $n_i > n_{i-1}$ such that $d(x, x_{n_i}) < 2^{-i}\delta$. Thus

$$d(q, x_{n_i}) < 2^{1-i}\delta$$

for $i = 1, 2, 3, \dots$. This says that (x_{n_i}) converges to q . Hence $q \in E$. \square

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. In the future you 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 8.15 (Cauchy sequence). A sequence (x_n) in X is a *Cauchy sequence* if

$$\forall \varepsilon > 0, \quad \exists N \in \mathbb{N}, \quad \forall n, m \geq N, \quad d(x_n, x_m) < \varepsilon.$$

Remark. Intuitively, we see that the distances between any two terms is sufficiently small after a certain point.

A natural question is regarding the relationship between convergent sequences and Cauchy sequences. We now address this.

Proposition 8.16.

- (i) In any metric space, every convergent sequence is a Cauchy sequence.
- (ii) If X is a compact metric space and if (x_n) is a Cauchy sequence in X , then (x_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 $x_n \rightarrow x$. Let $\varepsilon > 0$. There exists $N \in \mathbb{N}$ such that for all $n \geq N$,

$$d(x_n, x) < \frac{\varepsilon}{2}.$$

Then for all $n, m \geq N$,

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

as desired. Hence (x_n) is a Cauchy sequence.

- (ii) Let (x_n) be a Cauchy sequence in X . Since X is compact, it is sequentially compact. Then there exists a subsequence (x_{n_k}) such that $x_{n_k} \rightarrow x$.

Claim. $x_n \rightarrow x$.

Let $\varepsilon > 0$. Since (x_n) is a Cauchy sequence, there exists $N_1 \in \mathbb{N}$ such that

$$n, m \geq N_1 \implies d(x_n - x_m) < \frac{\varepsilon}{2}.$$

$x_{n_k} \rightarrow x$ implies there exists $N_2 \in \mathbb{N}$ such that

$$n_k \geq N_2 \implies d(x_{n_k}, x) < \frac{\varepsilon}{2}.$$

Let $N = \max\{N_1, N_2\}$, fix some $n_k \geq N$. Then $n \geq N$ implies

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

(iii) Suppose (x_n) is a Cauchy sequence.

We perform three steps:

- We first show that (x_n) is bounded:

Pick $N \in \mathbb{N}$ such that $|x_n - x_N| \leq 1$ for all $n \geq N$. Then

$$|x_n| \leq \max\{1 + |x_N|, |x_1|, \dots, |x_{N-1}|\}.$$

- Since (x_n) is bounded, by Bolzano–Weierstrass, (x_n) contains a subsequence (x_{n_k}) which converges to x .
- We now show that $x_n \rightarrow x$.

Let $\varepsilon > 0$ be given. Since (x_n) is a Cauchy sequence, there exists $N_1 \in \mathbb{N}$ such that

$$n, m \geq N_1 \implies |x_n - x_m| < \frac{\varepsilon}{2}.$$

Since $x_{n_k} \rightarrow x$, there exists $M \in \mathbb{N}$ such that for all $k > M$,

$$n_k > n_M \implies |x_{n_k} - x| < \frac{\varepsilon}{2}.$$

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

$$|x_n - x_{n_k}| < \frac{\varepsilon}{2}, \quad |x_{n_k} - x| < \frac{\varepsilon}{2}.$$

Hence

$$|x_n - x| \leq |x_n - x_{n_k}| + |x_{n_k} - x| < \varepsilon.$$

Therefore (x_n) is convergent, and $x_n \rightarrow x$.

□

Definition 8.17. A metric space X is *complete* if every Cauchy sequence in X converges.

Remark. The above result shows that 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 $x \in X$, and actually $x \in E$ since E is closed.)

Example 8.18. The sequence (x_n) is defined as follows:

$$x_n = 1 + \frac{1}{2} + \cdots + \frac{1}{n}.$$

(x_n) does not converge in \mathbb{R} .

Proof. We claim that (x_n) is not a Cauchy sequence. WLOG assume $n > m$. Consider

$$|x_n - x_m| = \frac{1}{m+1} + \frac{1}{m+2} + \cdots + \frac{1}{n} \geq \frac{n-m}{n} = 1 - \frac{m}{n}.$$

Let $n = 2m$, then

$$|x_n - x_m| = |x_{2m} - x_m| > \frac{1}{2}.$$

Hence (x_n) is not a Cauchy sequence, so it does not converge. □

Monotonic Sequences

Definition 8.19 (Monotonic sequence). A sequence (x_n) in \mathbb{R} is

- (i) *monotonically increasing* if $x_n \leq x_{n+1}$ for $n \in \mathbb{N}$;
- (ii) *monotonically decreasing* if $x_n \geq x_{n+1}$ for $n \in \mathbb{N}$;
- (iii) **monotonic** if it is either monotonically increasing or monotonically decreasing.

Lemma 8.20 (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 (x_n) is a monotonically increasing sequence bounded above. Let E be the range of x_n . Since E is bounded above, let $x = \sup E$.

Claim. $x_n \rightarrow x$.

By definition of supremum, $x_n \leq x$ for all $n \in \mathbb{N}$. For every $\varepsilon > 0$, there exists $N \in \mathbb{N}$ such that

$$x - \varepsilon < x_N \leq x,$$

otherwise $x - \varepsilon$ would be an upper bound of E . Since (x_n) is monotonically increasing, $n \geq N$ implies $x_N \leq x_n \leq x$, so

$$x - \varepsilon < x_n \leq x,$$

which implies $|x_n - x| < \varepsilon$. Hence $x_n \rightarrow x$. □

Limit Superior and Inferior

For divergent sequences, we have the following definition.

Definition 8.21. Suppose (x_n) is a sequence in \mathbb{R} . We write $x_n \rightarrow \infty$ if

$$\forall M \in \mathbb{R}, \quad \exists N \in \mathbb{N}, \quad \forall n \geq N, \quad x_n \geq M.$$

Similarly, we write $x_n \rightarrow -\infty$ if

$$\forall M \in \mathbb{R}, \quad \exists N \in \mathbb{N}, \quad \forall n \geq N, \quad x_n \leq M.$$

Definition 8.22. Suppose (x_n) is a sequence in \mathbb{R} . Let $E \subset \overline{\mathbb{R}}$ be the set of all subsequential limits of (x_n) (possibly including $+\infty$ and $-\infty$). Define

$$\limsup_{n \rightarrow \infty} x_n := \sup E,$$

$$\liminf_{n \rightarrow \infty} x_n := \inf E,$$

known as the *limit superior* and *limit infimum* of (x_n) respectively.

Remark. That is, limit superior is the “largest” subsequential limit; limit infimum is the “smallest” subsequential limit.

Remark. The limit superior and limit infimum exist due to the existence of supremum and infimum in $\overline{\mathbb{R}}$.

Equivalently, we can define limit superior (limit inferior) as the limit of supremum (infimum) of tails:

$$\begin{aligned} \limsup_{n \rightarrow \infty} x_n &= \lim_{n \rightarrow \infty} \left(\sup_{k \geq n} x_k \right) \\ &= \lim_{n \rightarrow \infty} (\sup \{a_n, a_{n+1}, \dots\}) \\ \liminf_{n \rightarrow \infty} x_n &= \lim_{n \rightarrow \infty} \left(\inf_{k \geq n} x_k \right) \\ &= \lim_{n \rightarrow \infty} (\inf \{a_n, a_{n+1}, \dots\}) \end{aligned}$$

Proposition 8.23. Suppose (x_n) is a sequence in \mathbb{R} . Then

(i) $\limsup_{n \rightarrow \infty} x_n \in E$;

(ii) if $x > \limsup_{n \rightarrow \infty} x_n$, there exists $N \in \mathbb{N}$ such that $x_n < x$ for all $n \geq N$.

Moreover, $\limsup_{n \rightarrow \infty} x_n$ is the only number that satisfies (i) and (ii).

Proof.

(i) We consider three cases for the value of $\limsup_{n \rightarrow \infty} x_n$:

- If $\limsup_{n \rightarrow \infty} x_n = +\infty$, then $\sup E = +\infty$, so E is not bounded above. Hence (x_n) is not bounded above, so (x_n) has a subsequence (x_{n_k}) such that $x_{n_k} \rightarrow \infty$.
- If $\limsup_{n \rightarrow \infty} x_n \in \mathbb{R}$, 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 $\limsup_{n \rightarrow \infty} x_n = -\infty$, then $\sup E = -\infty$, so E contains only one element, namely $-\infty$. Hence (x_n) has no subsequential limit. Thus for any $M \in \mathbb{R}$, $x_n > M$ for at most a finite number of values of n , so that $x_n \rightarrow -\infty$.

(ii) We prove by contradiction.

Suppose there is a number $x > \limsup_{n \rightarrow \infty} x_n$ such that $x_n \geq x$ for infinitely many values of n . In that case, there is a number $y \in E$ such that $y \geq x > \limsup_{n \rightarrow \infty} x_n$, contradicting the definition of $\limsup_{n \rightarrow \infty} x_n$.

We now show uniqueness. Suppose, for a contradiction, that two numbers p and q satisfy (i) and (ii). WLOG assume $p < q$. Then choose x such that $p < x < q$. Since p satisfies (i), we have $x_n < x$ for all $n \geq N$. But then q cannot satisfy (i). \square

Of course, an analogous result is true for $\liminf_{n \rightarrow \infty} x_n$.

Example 8.24. • Let (x_n) be a sequence containing all rationals. Then every real number is a subsequential limit, and

$$\limsup_{n \rightarrow \infty} x_n = +\infty, \quad \liminf_{n \rightarrow \infty} x_n = -\infty.$$

- Let $x_n = \frac{(-1)^n}{1 + \frac{1}{n}}$. Then

$$\limsup_{n \rightarrow \infty} x_n = 1, \quad \liminf_{n \rightarrow \infty} x_n = -1.$$

- For a sequence (x_n) in \mathbb{R} , $x_n \rightarrow x$ if and only if

$$\limsup_{n \rightarrow \infty} x_n = \liminf_{n \rightarrow \infty} x_n = x.$$

Proposition 8.25. If $a_n \leq b_n$ for $n \geq N$ where N is fixed, then

$$\liminf_{n \rightarrow \infty} a_n \leq \liminf_{n \rightarrow \infty} b_n,$$

$$\limsup_{n \rightarrow \infty} a_n \leq \limsup_{n \rightarrow \infty} b_n.$$

Proposition 8.26 (Arithmetic properties).

(i) If $k > 0$, $\limsup_{n \rightarrow \infty} ka_n = k \limsup_{n \rightarrow \infty} a_n$.

If $k < 0$, $\limsup_{n \rightarrow \infty} ka_n = k \liminf_{n \rightarrow \infty} a_n$.

$$(ii) \limsup(a_n + b_n) \leq \limsup a_n + \limsup b_n$$

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

§8.2 Series

Definition 8.27 (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 8.28 (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:

Proposition 8.29 (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 (Lemma 8.30)
- Boundedness of partial sums (Lemma 8.31, for series of non-negative terms)
- Comparison test (Lemma 8.32)
- Root test (Lemma 8.36)
- Ratio test (Lemma 8.37)
- Absolute convergence (Lemma 8.38)

Lemma 8.30 (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 8.31. *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 8.32 (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 8.33 (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 8.34 (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 $a_1 \geq a_2 \geq \cdots \geq 0$. Then $\sum a_n$ converges if and only if the series*

$$\sum_{k=0}^{\infty} 2^k a_{2^k} = a_1 + 2a_2 + 4a_4 + \cdots$$

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,

$$\begin{aligned} s_n &= a_1 + a_2 + \cdots + a_n, \\ t_k &= a_1 + 2a_2 + \cdots + 2^k a_{2^k}. \end{aligned}$$

We consider two cases:

- For $n < 2^k$, group terms to give

$$\begin{aligned} s_n &= a_1 + a_2 + \cdots + a_n \\ &\leq a_1 + (a_2 + a_3) + \cdots + (a_{2^k} + \cdots + a_{2^{k+1}-1}) \\ &\leq a_1 + 2a_2 + \cdots + 2^k a_{2^k} \\ &= t_k. \end{aligned}$$

By comparison test, if (t_k) converges, then (s_n) converges.

- For $n > 2^k$,

$$\begin{aligned} s_n &\geq a_1 + a_2 + (a_3 + a_4) + \cdots + (a_{2^{k-1}+1} + \cdots + a_{2^k}) \\ &\geq \frac{1}{2} a_1 + a_2 + 2a_4 + \cdots + 2^{k-1} a_{2^k} \\ &= \frac{1}{2} t_k. \end{aligned}$$

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 8.35 (The number e). Consider the series

$$\sum_{n=0}^{\infty} \frac{1}{n!}.$$

Claim. 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 $q!(e - s_n)$ 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 Proposition 8.25, 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 finally get

$$e \leq \liminf_{n \rightarrow \infty} t_n.$$

\square

Lemma 8.36 (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 8.37 (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 $n = N + k$,

$$|a_{N+1}| < \beta |a_N|$$

$$|a_{N+2}| < \beta |a_{N+1}| < \beta^2 |a_N|$$

\vdots

$$|a_{N+k}| < \beta^k |a_N|$$

Hence for all $n \geq N$,

$$\begin{aligned} |a_n| &< |a_N| \beta^{n-N} \\ &= (|a_N| \beta^{-N}) \beta^n \end{aligned}$$

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 of [Rud76].

The series $\sum a_n$ is said to *converge absolutely* if the series $\sum |a_n|$ converges.

Lemma 8.38 (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 8.39. 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 8.40 (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 8.41. *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

$$\begin{aligned} \left| \sum_{n=p}^q a_n b_n \right| &= \left| \sum_{n=p}^{q-1} A_n (b_n - b_{n+1}) + A_q b_q - A_{p-1} b_p \right| \\ &\leq M \left| \sum_{n=p}^{q-1} (b_n - b_{n+1}) + b_q + b_p \right| \quad [\cdot: |A_n| \leq M] \\ &= M |(b_p - b_q) + b_q + b_p| = 2M b_p \leq 2M b_n \leq \varepsilon. \end{aligned}$$

By the Cauchy criterion, $\sum a_n b_n$ converges to 0. □

Corollary 8.42 (Alternating series test). *Suppose (c_n) is a sequence such that*

- $|c_1| \geq |c_2| \geq |c_3| \geq \cdots$,
- $c_{2m-1} \geq 0, c_{2m} \leq 0$ 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;
- $b_0 \geq b_1 \geq b_2 \geq \cdots$ holds by assumption;
- $c_n \rightarrow 0$ implies $|c_n| \rightarrow 0$, so $b_n \rightarrow 0$.

Then by Proposition 8.41, we have that $\sum a_n b_n = 0$, so $\sum c_n = 0$. □

Addition and Multiplication of Series

Proposition 8.43. *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 8.44 (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 8.45 (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 8.46 (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 8.47 (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.

Proposition 8.48 (Riemann). *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. _____



to do

Proposition 8.49. *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.*

Exercises

Exercise 8.1. 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 the cases when $p > 1$, $p = 1$, and $0 < p < 1$.

If $p > 1$,

(iii)

(iv)

(v)

□

Exercise 8.2. Let (x_n) be a sequence in \mathbb{R} , 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 8.3 ([Rud76] 3.1). Prove that the convergence of (x_n) implies the convergence of $(|x_n|)$. Is the converse true?

Solution. Let $\varepsilon > 0$ be given. Since (x_n) is a Cauchy sequence, there exists $N \in \mathbb{N}$ such that for all $n, m \geq N$,

$$|x_n - x_m| < \varepsilon.$$

See that

$$||x_n| - |x_m|| \leq |x_n - x_m| < \varepsilon,$$

so $(|x_n|)$ is a Cauchy sequence, and therefore must converge.

The converse is not true, as shown by the sequence (x_n) with $x_n = (-1)^n$.

□

Exercise 8.4 ([Rud76] 3.3). The sequence (x_n) is recursively defined by

$$\begin{cases} x_0 = \sqrt{2}, \\ x_{n+1} = \sqrt{2 + x_n} \quad n \geq 0. \end{cases}$$

Show that (x_n) converges.

Solution. We first prove by induction that $x_n \leq x_{n+1} \leq 2$ for all $n \in \mathbb{N}$. For $n = 0$,

$$x_0 = \sqrt{2} \leq \sqrt{2 + \sqrt{2}} = x_1 \leq \sqrt{2 + \sqrt{4}} = 2.$$

If $x_{n-1} \leq x_n \leq 2$, then

$$x_n = \sqrt{2 + x_{n-1}} \leq \sqrt{2 + x_n} = x_{n+1} \leq \sqrt{2 + 2} = 2.$$

Hence (x_n) is monotonically increasing and bounded above by 2. By the monotone convergence theorem, (x_n) converges; let $x_n \rightarrow x$. Applying the limit on both sides of $x_{n+1} = \sqrt{2 + x_n}$,

$$\begin{aligned} \lim_{n \rightarrow \infty} x_{n+1} &= \lim_{n \rightarrow \infty} \sqrt{2 + x_n} \\ x &= \sqrt{2 + x} \\ x &= 2 \text{ or } 1 \end{aligned}$$

Since all $x_n \geq 0$, we must have $x = 2$. □

Exercise 8.5 (Contractive sequence). A sequence (x_n) in \mathbb{R} is *contractive* if there exists $k \in [0, 1)$ such that

$$|x_{n+2} - x_{n+1}| \leq k|x_{n+1} - x_n| \quad (\forall n \in \mathbb{N})$$

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 (in \mathbb{R}) and, hence, it is convergent. □

Exercise 8.6 ([Rud76] 3.4). Find the limit superior and limit inferior of the sequence (x_n) defined by

$$x_1 = 0, \quad x_{2m} = \frac{x_{2m-1}}{2}, \quad x_{2m+1} = x_{2m} + \frac{1}{2}.$$

Solution. We shall prove by induction that

$$x_{2m} = \frac{1}{2} - \frac{1}{2^m}, \quad x_{2m+1} = 1 - \frac{1}{2^m}$$

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 $x_2 = 0$ and $x_3 = \frac{1}{2}$. Hence assume that both formulae holds for $m \leq r$. Then

$$x_{2r+2} = \frac{1}{2}x_{2r+1} = \frac{1}{2} \left(1 - \frac{1}{2^r} \right) = \frac{1}{2} - \frac{1}{2^{r+1}}.$$

This completes the induction. We thus have $\limsup_{n \rightarrow \infty} x_n = 1$ and $\liminf_{n \rightarrow \infty} x_n = \frac{1}{2}$. □

Exercise 8.7 ([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 8.8 ([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 8.9 ([Rud76] 3.13). Prove that the Cauchy product of two absolutely convergent series converges absolutely.

Exercise 8.10 ([Rud76] 3.23). Suppose (p_n) and (q_n) are Cauchy sequences in a metric space X . Show that the sequence $(d(p_n, q_n))$ converges.

9 Continuity

§9.1 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 9.1 (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 9.2. 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 9.3. *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 Lemma 9.2, for every sequence (p_n) in $E \setminus \{p\}$ where $p_n \rightarrow p$, we have that

$$f(p_n) \rightarrow q \quad \text{and} \quad f(p_n) \rightarrow q'.$$

But the limit of a sequence is unique, so we must have $q = q'$. □

Lemma 9.4 (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 Lemma 9.2 and analogous properties of sequences in \mathbb{C} .

(i)

(ii)

(iii)

□

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 9.1, 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 9.5. For $c \in \mathbb{R}$, the set $\{x \in \mathbb{R} \mid x > c\}$ is called a neighbourhood of $+\infty$ and is written $(c, +\infty)$. Similarly, the set $(-\infty, c)$ is a neighbourhood of $-\infty$.

Definition 9.6. 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 $t \in V \cap E, t \neq x$.

to do

A moment's consideration will show that this coincides with Definition 4.1 when A and x are real.

The analogue of Theorem 4.4 is still true, and the proof offers nothing new. We state it, for the sake of completeness.

Lemma 9.7. Let $f, g : E \subset \mathbb{R} \rightarrow \mathbb{R}$. Suppose

§9.2 Continuous Functions

Definition 9.8 (Continuity). $f : E \subset X \rightarrow Y$ 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 9.1.)

Lemma 9.9. 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 9.1 and 9.8. □

If p is an isolated point of E , then our definition implies that every function f which has E as its domain of definition is continuous at p . For, no matter which $\varepsilon > 0$ we choose, we can pick $\delta > 0$ so that the only point $x \in E$ for which $d_X(x, p) < \delta$ is $x = p$; then $d_Y(f(x), f(p)) = 0 < \varepsilon$.

Corollary 9.10 (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 9.11. 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 Theorems 4.4 and 4.6 □

Example 9.12. 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 consider the composition of functions. The following result shows that a continuous function of a continuous function is continuous.

Proposition 9.13. 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.

Proposition 9.14. $f : X \rightarrow Y$ 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 $p \in f^{-1}(U)$. To show that $f^{-1}(U)$ is open in X , we will show that p is an interior point of $f^{-1}(U)$.

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 p is an interior point of $f^{-1}(U)$.

\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 9.15. $f : X \rightarrow Y$ 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

Definition 9.16. $f : E \rightarrow \mathbb{R}^k$ 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 9.17. 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 Proposition 9.14, 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 9.18. 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. From the previous result, $f(X)$ 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 9.19 (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 exists $p, q \in X$ 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 Proposition 7.32. □

Proposition 9.20. Suppose $f : X \rightarrow Y$ is continuous on X and bijective, X is compact. Then $f^{-1} : Y \rightarrow X$ is continuous on Y .

Proof. Applying Proposition 9.14 to f^{-1} in place of f , we see that to prove that f^{-1} is continuous on Y , it suffices to prove that $f(U)$ is open in Y for every open set U in X . Fix such a set U .

Since U is open in X , we have that U^c is closed in X . Since U^c is a closed subset of a compact set X , U^c is compact. Thus by Proposition 9.17, $f(U^c)$ 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 9.21 (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 9.22 (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 $c - \delta$ is an upper bound for A , again contradicting the definition of c .

Therefore we must have $f(c) = 0$. □

Continuity and Connectedness

Proposition 9.23. *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, i.e., separated. Then $A \cup B = f(E)$ 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 Corollary 9.15, $f^{-1}(\overline{A})$ 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 9.24 (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 Proposition 7.66, $[a, b]$ is connected. By the previous result, we have that $f([a, b])$ is a connected subset of \mathbb{R} . Then apply Proposition 7.67 and we are done. □

Remark. The converse is not necessarily true. For instance, consider the *topologist's sine curve*:

$$f(x) = \begin{cases} 0 & (x = 0) \\ \sin\left(\frac{1}{x}\right) & (x \neq 0) \end{cases}$$

f satisfies the intermediate value property but f is not continuous.

§9.3 Uniform Continuity

Definition 9.25 (Uniform continuity). $f : X \rightarrow Y$ 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 9.26. • 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 $|p - q| < \delta$. This shows that f is uniformly continuous on $(0, 1]$.

The next result concerns the relationship between continuity and uniform continuity.

Lemma 9.27.

- (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 9.28 (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

§9.4 Discontinuities

Definition 9.29 (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 9.1; 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 9.30. 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 9.31 (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 9.32 (Discontinuities). Let $f : (a, b) \rightarrow \mathbb{R}$. Suppose f is discontinuous at $x \in (a, b)$.

- (i) f has a *discontinuity of the first kind* (or a *simple discontinuity*) at x , if $f(x+)$ and $f(x-)$ exist;
- (ii) otherwise f has a *discontinuity of the second kind*.

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 9.33. • 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.
- 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 $(-3, 1)$.

§9.5 Monotonic Functions

We now study those functions which never decrease (or never increase) on a given interval.

Definition 9.34 (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 9.35. 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 9.36. *Monotonic functions have no discontinuities of the second kind.*

Proposition 9.37. *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. □

§9.6 Lipschitz Continuity

Definition 9.38. $f : X \rightarrow Y$ 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 ; f may also be referred to as *K -Lipschitz*.

Lemma 9.39. *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.

- First, suppose that $K \leq 0$. 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.

- Next, suppose that $K > 0$. Take $\delta = \frac{\varepsilon}{K}$. If $d_X(x, y) < \delta$, then

$$K d_X(x, y) < \varepsilon.$$

By Lipschitz continuity of f , we have that

$$d_Y(f(x), f(y)) \leq K d_X(x, y).$$

These last two statements together imply $d_Y(f(x), f(y)) < \varepsilon$. Hence f is uniformly continuous on X .

□

$f : X \rightarrow Y$ is a *contraction* if it is a K -Lipschitz map for some $K < 1$.

If $f : X \rightarrow X$ is a map, $x \in X$ is called a *fixed point* if $f(x) = x$.

Theorem 9.40 (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, the sequence (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 y also be a fixed point. Then

$$d(x, y) = d(f(x), f(y)) = Kd(x, y).$$

As $K < 1$ this means that $d(x, y) = 0$ and hence $x = y$. 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.

Exercises

Exercise 9.1 ([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 9.2 ([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 9.3 ([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 9.4 ([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 9.5 ([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 9.6 ([Rud76] 4.12). A uniformly continuous function of a uniformly continuous function is uniformly continuous.

Exercise 9.7 ([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 9.8 ([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 9.9 ([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 9.10 ([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 9.11 ([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.

10 Differentiation

§10.1 The Derivative of A Real Function

Definition 10.1 (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 10.2 (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 Lemma 9.9, f is continuous at x . □

Remark. The converse of Lemma 10.2 is not true; it is easy to construct continuous functions which fail to be differentiable at isolated points.

Example 10.3 (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 10.4. 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.

Notation. 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. $\mathcal{C}_1[a, b]$ denotes the set of differentiable functions over $[a, b]$ whose derivative is continuous. More generally, $\mathcal{C}_n[a, b]$ denotes the set of functions whose n -th derivative is continuous. In particular, $\mathcal{C}_0[a, b]$ is the set of continuous functions over $[a, b]$.

Lemma 10.5 (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) $\frac{f}{g}$ (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 10.6. 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 10.7 (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*

$$h'(x) = g'(f(x)) f'(x). \quad (10.1)$$

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) .

§10.2 Mean Value Theorems

Let (X, d) be a metric space.

Definition 10.8 (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 10.9 (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 $f'(x) = 0$. □

Theorem 10.10 (Rolle's theorem). Suppose f is continuous on $[a, b]$, 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 (Theorem 9.19), 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 10.11 (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

$$[f(b) - f(a)]g'(c) = [g(b) - g(a)]f'(c). \quad (10.2)$$

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 10.12 (Mean value theorem). Suppose f is continuous on $[a, b]$ and differentiable in (a, b) . Then there exists $c \in (a, b)$ such that

$$f(b) - f(a) = f'(c)(b - a). \quad (10.3)$$

Proof. Take $g(x) = x$ in Theorem 10.11. □

Lemma 10.13. 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 . □

§10.3 Continuity of Derivatives

The following result implies some sort of a “intermediate value” property of derivatives that is similar to continuous functions.

Theorem 10.14 (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 10.15. *If f is differentiable on $[a, b]$, then f' cannot have any simple discontinuities on $[a, b]$.*

§10.4 L'Hopital's Rule

The following result is frequently used in the evaluation of limits.

to do

Lemma 10.16 (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

$$\lim_{x \rightarrow a} \frac{f(x)}{g(x)} = A.$$

Proof. The entire proof is rather tedious because we have too many cases.

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 $A < r < q$.

1. $\frac{0}{0}$ or $\frac{\infty}{\infty}$ 2. a is normal or $a = -\infty$ 3. A is normal or $A = \pm\infty$

We'll only prove the most basic one here: $0/0$, a and A are normal This is the case which will be required for Taylor series

First we define $f(a)=g(a)=0$, so that f and g are continuous at $x = a$

Now let $x \in (a, b)$, then f and g are continuous on $[a, x]$ and differentiable in (a, x) : Thus by Cauchy's Mean Value Theorem, there exists $\xi \in (a, x)$ such that

$$\frac{f'(\xi)}{g'(\xi)} = \frac{f(x) - f(a)}{g(x) - g(a)} = \frac{f(x)}{g(x)}$$

For each x , we pick ξ which satisfies the above, so that ξ may be seen as a function of x satisfying $a < \xi(x) < x$

Then by squeezing we have $\lim_{x \rightarrow a^+} \xi(x) = a$.

Since $\frac{f'}{g'}$ is continuous near a , the theorem regarding the limit of composite functions give

$$\lim_{x \rightarrow a^+} \frac{f(x)}{g(x)} = \lim_{x \rightarrow a^+} \frac{f'(\xi)}{g'(\xi)} = \lim_{x \rightarrow a^+} \left(\frac{f'}{g'} \right) (\xi(x)) = A$$

Now the same reasoning can be used for b where we will use $\lim_{x \rightarrow b^-}$ to replace all the $\lim_{x \rightarrow a^+}$, and ξ will be a function which maps to (x, b) . \square

§10.5 Taylor's Theorem

Theorem 10.17 (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

$$f(x) = P_{n-1}(x) + \frac{f^{(n)}(z_x)}{n!} (x - c)^n. \quad (10.4)$$

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. (10.4) 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 10.18.

$$\begin{aligned}e^x &= 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \cdots \\ \sin x &= x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \cdots \\ \cos x &= 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \cdots \\ \ln(1+x) &= x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \cdots\end{aligned}$$

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).

§10.6 Differentiation of Vector-valued Functions

Let $\mathbf{f} : (a, b) \rightarrow \mathbb{R}^n$. Then $\mathbf{f} = (f_1, \dots, f_n)$, where each component $f_k : (a, b) \rightarrow \mathbb{R}$. We say that \mathbf{f} is differentiable at $x \in (a, b)$ if each component f_k is differentiable at x :

$$\mathbf{f}'(x) = (f'_1(x), \dots, f'_n(x)).$$

Exercises

Exercise 10.1. 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 10.2. Given that $f(x) = x^\alpha$ where $0 < \alpha < 1$. Prove that f is uniformly continuous on $[0, +\infty)$.

Exercise 10.3. 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

$$f(x) + f'(x) = 0.$$

11 Riemann–Stieltjes Integral

§11.1 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, say

$$P = \{x_0, x_1, \dots, x_n\},$$

where $a = x_0 \leq x_1 \leq \dots \leq x_{n-1} \leq x_n = b$.

Notation. $\mathcal{P}[a, b]$ denotes the set of all partitions of $[a, b]$.

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

$$\Delta\alpha_i = \alpha(x_i) - \alpha(x_{i-1}) \quad (i = 1, \dots, n).$$

The *upper sum* and *lower sum* of f with respect to the partition P and α are respectively

$$U(f, \alpha; P) = \sum_{i=1}^n M_i \Delta\alpha_i,$$

$$L(f, \alpha; P) = \sum_{i=1}^n m_i \Delta\alpha_i.$$

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. We will now show this.

Lemma 11.1. *If P' is a refinement of P , then*

$$(i) \quad L(f, \alpha; P) \leq L(f, \alpha; P')$$

$$(ii) \quad U(f, \alpha; P') \leq U(f, \alpha; P)$$

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

$$\begin{aligned} \int_a^b f \, d\alpha &:= \inf_{P \in \mathcal{P}[a, b]} U(f, \alpha; P) \\ \int_a^b f \, d\alpha &:= \sup_{P \in \mathcal{P}[a, b]} L(f, \alpha; P) \end{aligned}$$

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 11.2.

$$\int_a^b f \, d\alpha \leq \int_a^{\bar{b}} f \, d\alpha .$$

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 Lemma 11.1,

$$L(f, \alpha; P_1) \leq L(f, \alpha; P').$$

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

$$\int_a^b f \, d\alpha \leq U(f, \alpha; P_2).$$

Then taking inf over all P_2 gives

$$\int_a^b f \, d\alpha \leq \int_a^{\bar{b}} f \, d\alpha .$$

□

Definition

Definition 11.3 (Riemann–Stieltjes integral). f 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. $\mathcal{R}_\alpha[a, b]$ denotes the set of Riemann–Stieltjes integrable functions with respect to α over $[a, b]$.

In particular, when $\alpha(x) = x$, we call the corresponding Riemann–Stieltjes integration the *Riemann integration*, and use $\mathcal{R}[a, b]$ to denote the set of Riemann integrable functions.

Notation. x is a “dummy variable” and may be replaced by any other convenient symbol; hence we omit it.

Example 11.4 (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 $m_i = 0$. Hence for any partition P ,

$$U(f; P) = 1, \quad L(f; P) = 0.$$

Thus

$$1 = \int_a^b f \, d\alpha \neq \int_a^b f \, d\alpha = 0$$

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 11.5 (Integrability criterion). $f \in \mathcal{R}_\alpha[a, b]$ if and only if

$$\forall \varepsilon > 0, \quad \exists P, \quad U(f, \alpha; P) - L(f, \alpha; P) < \varepsilon.$$

Proof.

\Rightarrow Suppose $f \in \mathcal{R}_\alpha[a, b]$. 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

$$U(f, \alpha; P) - L(f, \alpha; P) < \varepsilon.$$

\Leftarrow From Lemma 11.2, for every partition P , we have

$$L(f, \alpha; P) \leq \int_a^b f \, d\alpha \leq \int_a^b f \, d\alpha \leq U(f, \alpha; P).$$

Since $U(f, \alpha; P) - L(f, \alpha; P) < \varepsilon$, we have that

$$0 \leq \int_a^{\bar{b}} d\alpha - \int_a^b f d\alpha < \varepsilon.$$

Since this holds for all $\varepsilon > 0$, we have

$$\int_a^{\bar{b}} f d\alpha = \int_a^b f d\alpha.$$

Hence $f \in \mathcal{R}_\alpha[a, b]$. □

Useful Identities

Proposition 11.6 (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 $f \in \mathcal{R}_\alpha[a, b]$ 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 Lemma 11.1, 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) See that

$$f(s_i), f(t_i) \in [m_i, M_i] \quad (i = 1, \dots, n)$$

so that

$$|f(s_i) - f(t_i)| \leq M_i - m_i.$$

Thus

$$\sum_{i=1}^n |f(s_i) - f(t_i)| \Delta\alpha_i \leq U(f, \alpha; P) - L(f, \alpha; P) < \varepsilon.$$

(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 11.7 (Continuity implies integrability). *If f is continuous on $[a, b]$, then $f \in \mathcal{R}_\alpha[a, b]$.*

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 Lemma 9.27, 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 $f \in \mathcal{R}_\alpha[a, b]$, by the integrability criterion (Lemma 11.5). □

Proposition 11.8. *If f is monotonic on $[a, b]$, and if α is continuous on $[a, b]$, then $f \in \mathcal{R}_\alpha[a, b]$.*

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 $f \in \mathcal{R}_\alpha[a, b]$, by the integrability criterion. \square

Proposition 11.9. 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 $f \in \mathcal{R}_\alpha[a, b]$.

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 $f \in \mathcal{R}_\alpha[a, b]$, by the integrability criterion. \square

The next result states that a uniformly continuous function of an integrable function is also integrable.

Proposition 11.10. Suppose $f \in \mathcal{R}_\alpha[a, b]$, $m \leq f \leq M$, and ϕ is continuous on $[m, M]$. Then $\phi \circ f \in \mathcal{R}_\alpha[a, b]$.

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 $f \in \mathcal{R}_\alpha[a, b]$, by Lemma 11.5, there exists a partition $P = \{x_0, \dots, x_n\}$ of $[a, b]$ such that

$$U(f, \alpha; P) - L(f, \alpha; P) < \delta^2. \quad (1)$$

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, $h \in \mathcal{R}_\alpha[a, b]$. □

§11.2 Properties of the Integral

Lemma 11.11.

(i) If $f_1, f_2 \in \mathcal{R}_\alpha[a, b]$, then $f_1 + f_2 \in \mathcal{R}_\alpha[a, b]$, 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 $f \in \mathcal{R}_\alpha[a, b]$, then $cf \in \mathcal{R}_\alpha[a, b]$ for every $c \in \mathbb{R}$, and

$$\int_a^b (cf) d\alpha = c \int_a^b f d\alpha.$$

(iii) If $f_1, f_2 \in \mathcal{R}_\alpha[a, b]$ and $f_1 \leq f_2$, then

$$\int_a^b f_1 d\alpha \leq \int_a^b f_2 d\alpha.$$

(iv) If $f \in \mathcal{R}_\alpha[a, b]$ 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 $f \in \mathcal{R}_\alpha[a, b]$ 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 $f \in \mathcal{R}_\alpha[a, b]$ 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 $f \in \mathcal{R}_\alpha[a, b]$ and $g \in \mathcal{R}_\alpha[a, b]$, then $fg \in \mathcal{R}_\alpha[a, b]$.

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 $f_1 \in \mathcal{R}_\alpha[a, b]$ and $f_2 \in \mathcal{R}_\alpha[a, b]$, 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 $f \in \mathcal{R}_\alpha[a, b]$.

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 $cf \in \mathcal{R}_\alpha[a, b]$.

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 Proposition 11.10, $f^2 \in \mathcal{R}_\alpha[a, b]$ if $f \in \mathcal{R}_\alpha[a, b]$. Write

$$fg = \frac{1}{4} \left((f+g)^2 - (f-g)^2 \right).$$

Then the desired result follows. □

Lemma 11.12 (Triangle inequality). Suppose $f \in \mathcal{R}_\alpha[a, b]$. Then $|f| \in \mathcal{R}_\alpha[a, b]$, 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 Proposition 11.10, we have that $|f| = \phi \circ f \in \mathcal{R}_\alpha[a, b]$. 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 11.13 (Heaviside step function). The Heaviside step function H 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 $x_2 \rightarrow s$. \square

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 Theorems 6.12 and 6.15,

$$\int_a^b f \, d\alpha_1 = \sum_{n=1}^N c_n f(s_n).$$

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 $\alpha = \alpha_1 + \alpha_2$, it follows from (24) and (25) 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

$$\int_a^b f \, d\alpha = \sum_{n=1}^{\infty} c_n f(s_n).$$

\square

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 11.14 (Integration by substitution). *Assume α increases monotonically, $\alpha' \in \mathcal{R}[a, b]$. Let $f : [a, b] \rightarrow \mathbb{R}$ be bounded, then $f \in \mathcal{R}_\alpha[a, b]$ if and only if $f\alpha' \in \mathcal{R}[a, b]$. In that case*

$$\int_a^b f \, d\alpha = \int_a^b f(x)\alpha'(x) \, dx. \quad (11.1)$$

Proof. Let $\varepsilon > 0$ be given and apply Theorem 6.6 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 Proposition 11.6,

$$\sum_{i=1}^n \left| \alpha'(s_i) - \alpha'(t_i) \right| \Delta x_i < \varepsilon. \quad (2)$$

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

$$\begin{aligned} \left| \sum_{i=1}^n f(s_i)\Delta\alpha_i - \sum_{i=1}^n f(s_i)\alpha'(s_i)\Delta x_i \right| &= \left| \sum_{i=1}^n f(s_i) \left(\alpha'(t_i) - \alpha'(s_i) \right) \Delta x_i \right| \\ &\leq \sum_{i=1}^n \left| f(s_i) \left(\alpha'(t_i) - \alpha'(s_i) \right) \Delta x_i \right| \\ &= \sum_{i=1}^n |f(s_i)| \left| \alpha'(t_i) - \alpha'(s_i) \right| \Delta x_i \\ &\leq M \sum_{i=1}^n \left| \alpha'(t_i) - \alpha'(s_i) \right| \Delta x_i \\ &\leq M\varepsilon. \end{aligned} \quad (3)$$

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

$$\left| U(f, \alpha; P) - U(f\alpha'; P) \right| \leq M\varepsilon. \quad (4)$$

Since (1) holds true for any refinement of P , hence (4) also remains true. We conclude that

$$\left| \int_a^b f \, d\alpha - \int_a^b f(x)\alpha'(x) \, dx \right| \leq M\varepsilon.$$

But ε is arbitrary. Hence

$$\int_a^b f \, d\alpha = \int_a^b f(x)\alpha'(x) \, dx$$

for any bounded f . The equality of the lower integrals follows from

$$\begin{aligned} \int_a^b -f \, d\alpha &= \int_a^b -f\alpha' \, dx \\ - \int_a^b f \, d\alpha &= - \int_a^b f\alpha' \, dx \\ \int_a^b f \, d\alpha &= \int_a^b f(x)\alpha'(x) \, dx \end{aligned}$$

Therefore the theorem follows. □

Proposition 11.15 (Change of variables). *Suppose $\phi : [A, B] \rightarrow [a, b]$ is strictly increasing and continuous. Suppose α is monotonically increasing on $[a, b]$, $f \in \mathcal{R}_\alpha[a, b]$. Define β and g on $[A, B]$ by*

$$\beta(y) = \alpha(\phi(y)), \quad g(y) = f(\phi(y)).$$

Then $g \in \mathcal{R}_\beta[A, B]$, and

$$\int_A^B g \, d\beta = \int_a^b f \, d\alpha. \quad (11.2)$$

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

$$\begin{aligned} U(g, \beta; Q) &= U(f, \alpha; P), \\ L(g, \beta; Q) &= L(f, \alpha; P). \end{aligned}$$

Since $f \in \mathcal{R}_\alpha[a, b]$, P can be chosen so that both $U(f, \alpha; P)$ and $L(f, \alpha; P)$ are close to $\int f \, d\alpha$. Hence (38), combined with Theorem 6.6, shows that $g \in \mathcal{R}_\beta[A, B]$ and that (37) holds. This completes the proof. □

Note the following special case: Take $\alpha(x) = x$. Then $\beta = \phi$. Assume $\phi' \in \mathcal{R}[A, B]$. If Theorem 6.17 is applied to the left side of (37), we obtain

$$\int_a^b f(x) \, dx = \int_A^B f(\phi(y)) \phi'(y) \, dy.$$

§11.3 Integration and Differentiation

We shall show that integration and differentiation are, in a certain sense, inverse operations.

Theorem 11.16. Suppose $f \in \mathcal{R}_\alpha[a, b]$. 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 $f \in \mathcal{R}_\alpha[a, b]$. 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 Theorem 6.12(d),

$$\begin{aligned} \left| \frac{F(t) - F(s)}{t - s} - f(x_0) \right| &= \left| \frac{\int_a^s f(u) \, du - \int_a^s f(u) \, du}{t - s} - f(x_0) \right| \\ &= \left| \frac{1}{t - s} \int_s^t (f(u) - f(x_0)) \, du \right| \\ &= \frac{1}{t - s} \left| \int_s^t (f(u) - f(x_0)) \, du \right| \\ &\leq \frac{1}{t - s} \int_s^t |f(u) - f(x_0)| \, du \\ &< \frac{1}{t - s} \varepsilon(t - s) = \varepsilon \end{aligned}$$

so it follows that $F'(x_0) = f(x_0)$. □

Theorem 11.17 (Fundamental theorem of calculus). *Suppose $f \in \mathcal{R}_\alpha[a, b]$, and there exists a differentiable function F on $[a, b]$ such that $F' = f$. Then*

$$\int_a^b f(x) \, dx = F(b) - F(a). \quad (11.3)$$

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 Proposition 11.6,

$$\left| F(b) - F(a) - \int_a^b f(x) \, dx \right| = \left| \sum_{i=1}^n f(t_i) \Delta x_i - \int_a^b f(x) \, dx \right| < \varepsilon.$$

Since this holds for all $\varepsilon > 0$, the proof is complete. \square

Lemma 11.18 (Integration by parts). *Suppose F and G are differentiable on $[a, b]$, $F' = f \in \mathcal{R}[a, b]$ and $G' = g \in \mathcal{R}[a, b]$. Then*

$$\int_a^b F(x)g(x) \, dx = F(b)G(b) - F(a)G(a) - \int_a^b f(x)G(x) \, dx. \quad (11.4)$$

Proof. Let $H(x) = F(x)G(x)$. Then

$$\begin{aligned} H'(x) &= F'(x)G(x) + F(x)G'(x) \\ &= f(x)G(x) + F(x)g(x) \end{aligned}$$

\square

§11.4 Integration of Vector-valued Functions

Definition 11.19. Let $f_1, \dots, f_k : [a, b] \rightarrow \mathbb{R}$, and let $\mathbf{f} = (f_1, \dots, f_k) : [a, b] \rightarrow \mathbb{R}^k$. We say that $\mathbf{f} \in \mathcal{R}_\alpha[a, b]$ if $f_1, \dots, f_k \in \mathcal{R}_\alpha[a, b]$. 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, $\int \mathbf{f} \, d\alpha$ 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.

Proposition 11.20. If $\mathbf{f}, \mathbf{F} : [a, b] \rightarrow \mathbb{R}^k$, $\mathbf{f} \in \mathcal{R}_\alpha[a, b]$, and $\mathbf{F}' = \mathbf{f}$. Then

$$\int_a^b \mathbf{f}(t) \, dt = \mathbf{F}(b) - \mathbf{F}(a). \quad (11.5)$$

to do

§11.5 Rectifiable Curves

Definition 11.21 (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 $[a, b]$, let

$$\Lambda(\gamma; P) = \sum_{i=1}^n |\gamma(x_i) - \gamma(x_{i-1})|.$$

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 11.22. The *total variation* (or *length*) of γ is

$$\Lambda(\gamma) := \sup_{P \in \mathcal{P}[a, b]} \Lambda(\gamma; P).$$

Definition 11.23 (Rectifiable curve). γ is *rectifiable* if $\Lambda(\gamma) < \infty$.

Proposition 11.24. If γ is a continuously differentiable curve on $[a, b]$, then γ is rectifiable, and

$$\Lambda(\gamma) = \int_a^b |\gamma'(t)| dt. \quad (11.6)$$

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

$$\Lambda(\gamma) \leq \int_a^b |\gamma'(t)| dt.$$

To prove the opposite inequality, let $\varepsilon > 0$ be given. Since γ' is uniformly continuous on $[a, b]$, 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,

$$\int_a^b |\gamma'(t)| \, dt \leq \Lambda(\gamma).$$

This completes the proof. □

Exercises

12 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 $Y = \mathbb{C}$).

§12.1 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 12.1 (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 12.2. 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 12.3 (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 12.4 (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 12.5 (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.

§12.2 Uniform Convergence

Definition 12.6 (Uniform convergence). (f_n) *converges uniformly* to f over E , denoted by $f_n \Rightarrow f$, 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 .

Remark. 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 .

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 .

Remark. Uniform convergence implies pointwise convergence, but not the other way around.

Example 12.7. Consider the sequence of functions $f_n(x) = x^n$ defined on $(0, 1)$. Then $f_n \rightarrow 0$. But $f_n \not\Rightarrow 0$.

Proof. □

The Cauchy criterion for uniform convergence is as follows.

Lemma 12.8 (Cauchy criterion). Suppose (f_n) is a sequence of complex-valued functions. Then $f_n \Rightarrow f$ on E if and only if (f_n) is uniformly Cauchy:

$$\forall \varepsilon > 0, \quad \exists N \in \mathbb{N}, \quad \forall x \in E, \quad \forall n, m \geq N, \quad |f_n(x) - f_m(x)| < \varepsilon.$$

Proof.

\Rightarrow Suppose $f_n \Rightarrow f$ 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 (f_n) is uniformly Cauchy. 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. \square

Definition 12.9. Let X and Y be metric spaces. Then $\mathcal{C}(X, Y)$ denotes the space of continuous bounded functions from X to Y . If $f \in \mathcal{C}(X, \mathbb{C})$, we define the **supremum norm** of f as

$$\|f\| := \sup_{x \in X} |f(x)|.$$

Lemma 12.10. $\|f\|$ 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. \square

The following criterion is sometimes useful.

Lemma 12.11. Suppose (f_n) is a sequence of complex-valued functions defined on E . Then $f_n \Rightarrow f$ 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 $f_n \Rightarrow f$ on E , by definition.

Note that for the last step, the \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 12.12 (Weierstrass M-test). *Suppose (f_n) is a sequence of 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 (Lemma 12.8), we are done. \square

Example 12.13. • 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 $[-r, r] \subset \mathbb{R}$,

$$\left| \frac{x^n}{n!} \right| \leq \frac{r^n}{n!} \quad \text{and} \quad \sum_{n=1}^{\infty} \frac{r^n}{n!} \text{ converges by the ratio test.}$$

§12.3 Properties of Uniform Convergence

We now consider properties preserved by uniform convergence.

Uniform Convergence and Continuity

We prove a more general result.

Proposition 12.14. *Suppose (f_n) is a sequence of complex-valued functions defined on E , such that $f_n \Rightarrow f$ 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. We first show that (A_n) converges. Since (f_n) uniformly converges on E , (f_n) is uniformly Cauchy. Let $\varepsilon > 0$ be given, 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 .

Next we will show that $\lim_{t \rightarrow x} f(t) = A$. 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 $\lim_{t \rightarrow x} f(t) = A$. \square

An immediate important corollary is that uniform convergence preserves continuity.

Corollary 12.15. *Suppose (f_n) are continuous on E , and $f_n \Rightarrow f$ 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 . \square

Remark. The converse is not true. Just because the limit is continuous doesn't mean that the convergence is uniform. For example, $f_n : (0, 1) \rightarrow \mathbb{R}$ defined by $f_n(x) = x^n$ converges to the zero function, but not uniformly.

Let us see that we can have extra conditions such that the converse is true.

Proposition 12.16 (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 $f_n \Rightarrow f$ on K .

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 $g_n \Rightarrow 0$ 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 the set $\{g_n(x) \mid g_n(x) \geq \varepsilon\}$ is closed, by Corollary 9.15, its pre-image K_n is closed. Since K_n is a closed subset of a compact set K , then K_n is compact (by Proposition 7.40).

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 $0 \leq g_n(x) < \varepsilon$ for all $x \in K$ and for all $n \geq N$. This proves the theorem. \square

Remark. The compactness in the hypotheses is necessary; for instance, for $0 < x < 1$, let

$$f_n(x) = \frac{1}{nx + 1} \quad (n = 1, 2, \dots).$$

Then $f_n(x) \rightarrow 0$ monotonically in $(0, 1)$, but the convergence is not uniform.

Lemma 12.17. $\mathcal{C}(X, \mathbb{C})$ 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 (Lemma 12.8), $f_n \Rightarrow f$ 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 Corollary 12.15.
- 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 $f_n \Rightarrow f$ on X , we have $\|f - f_n\| \rightarrow 0$ as $n \rightarrow \infty$. □

Uniform Convergence and Integration

The next result states that the limit and integral can be interchanged.

Proposition 12.18. *Suppose (f_n) is a sequence of functions defined over $[a, b]$ and $f_n \in \mathcal{R}_\alpha[a, b]$. If $f_n \Rightarrow f$ on $[a, b]$, then $f \in \mathcal{R}_\alpha[a, b]$, 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

$$\varepsilon_n = \sup_{x \in [a, b]} |f_n(x) - f(x)|.$$

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

$$0 \leq \int_a^b f \, d\alpha - \int_a^b f_n \, d\alpha \leq 2\varepsilon_n[\alpha(b) - \alpha(a)].$$

Since $\varepsilon_n \rightarrow 0$ as $n \rightarrow \infty$ (Theorem 7.9), the upper and lower integrals of f are equal. Thus $f \in \mathcal{R}_\alpha[a, b]$.

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 12.19. *Suppose $f_n \in \mathcal{R}_\alpha[a, b]$ 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 $f_n \in \mathcal{R}_\alpha[a, b]$ and $f_n \Rightarrow f$. Apply above theorem to (f_n) and the conclusion follows. \square

Example 12.20. 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

Proposition 12.21. 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 $f_n \Rightarrow f$ 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 $[a, b]$, then (f'_n) is uniformly Cauchy (by Lemma 12.8), so

$$|f'_n(x) - f'_m(x)| < \frac{\varepsilon}{2(b-a)} \quad (a \leq x \leq b).$$

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

$$\begin{aligned} |(f_n(x) - f_m(x)) - (f_n(x_0) - f_m(x_0))| &= |f'_n(t) - f'_m(t)| |x_0 - x| \\ &< \frac{\varepsilon}{2(b-a)} |x_0 - x| \\ &\leq \frac{\varepsilon}{2} \end{aligned} \tag{1}$$

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 Lemma 12.8, (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 12.22 (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 .

§12.4 Equicontinuous Families of Functions

We would like an analogue of Bolzano–Weierstrass; that is, every bounded sequence of functions has a convergent subsequence.

Definition 12.23. 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 12.24. 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 (c), the sequence S (except possibly its first $n - 1$ terms) is a subsequence of S_n , for $n = 1, 2, \dots$. Hence (b) implies that $(f_{n,n}(x_i))$ converges, as $n \rightarrow \infty$, for every $x_i \in E$. □

to do

Definition 12.25. 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 12.26. 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 $d(x, y) < \delta$, it follows that

$$\begin{aligned} |f_n(x) - f_n(y)| &\leq |f_n(x) - f_N(x)| + |f_N(x) - f_N(y)| + |f_N(y) - f_n(y)| \\ &< \frac{\varepsilon}{3} + \frac{\varepsilon}{3} + \frac{\varepsilon}{3} = \varepsilon. \end{aligned}$$

In conjunction with (43), this proves the theorem. □

We first need the following lemma.

Lemma 12.27. 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 12.28 (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 12.29. *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 12.30. *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}) .*

§12.5 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 12.31 (Weierstrass approximation theorem). *If $f : [a, b] \rightarrow \mathbb{C}$ is continuous, there exists a sequence of polynomials (P_n) such that $P_n \Rightarrow f$ 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 12.32. *The metric space $\mathcal{C}([a, b], \mathbb{C})$ contains a countable dense subset.*

Corollary 12.33. *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 12.34. 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 12.35.

Proposition 12.36. 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 12.37. 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 12.38.

Proposition 12.39. 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 12.40 (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

13 Some Special Functions

§13.1 Power Series

Definition 13.1. 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 13.2 (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

$$D_R(a) := \{z \in \mathbb{C} : |z| < R\}.$$

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.

(ii) If $|z| > R$, then $\limsup_{n \rightarrow \infty} \sqrt[n]{|a_n|} > 1$. By the root test, $\sum c_n z^n$ diverges.

□

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 13.3. 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 13.4. 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 8.41 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 13.5. 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 13.6. 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 uniformly on $[-R + \varepsilon, R - \varepsilon]$ for all $\varepsilon > 0$;
- (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) Let $\varepsilon > 0$ be given. For $|x| \leq R - \varepsilon$, we have

$$|c_n x^n| \leq |c_n (R - \varepsilon)^n|$$

and since

$$\sum c_n (R - \varepsilon)^n$$

converges absolutely (every power series converges absolutely in the interior of its interval of convergence, by the root test), Theorem 7.10 show that $\sum c_n x^n$ uniformly converges on $[-R + \varepsilon, R - \varepsilon]$.

- (ii) Since $\sqrt[n]{n} \rightarrow 1$ as $n \rightarrow \infty$, we have

$$\limsup_{n \rightarrow \infty} \sqrt[n]{n|c_n|} = \limsup_{n \rightarrow \infty} \sqrt[n]{|c_n|},$$

so that the series (4) and (5) have the same interval of convergence. Since (5) is a power series, it converges uniformly in $[-R + \varepsilon, R - \varepsilon]$, for every $\varepsilon > 0$, and we can apply Theorem 7.17 (for series instead of sequences). It follows that (5) holds if $|x| \leq R - \varepsilon$.

But, given any x such that $|x| < R$, we can find an $\varepsilon > 0$ such that $|x| < R - \varepsilon$. This shows that (5) holds for $|x| < R$.

Continuity of f follows from the differentiability of f on $(-R, R)$.

□

Corollary 13.7. f is infinitely differentiable in $(-R, R)$; its derivatives are given by

$$f^{(k)}(x) = \sum_{n=k}^{\infty} n(n-1) \cdots (n-k+1) c_n x^{n-k}.$$

In particular,

$$f^{(k)}(0) = k! c_k, \quad k = 0, 1, 2, \dots$$

Proof. Apply theorem successively to f, f', f'', \dots . Put $x = 0$.

□

If the series (3) converges at an endpoint, say at $x = R$, then/is continuous not only in $(-R, R)$, but also at $x = R$. This follows from the following result (for simplicity of notation, we take $R = 1$).

Proposition 13.8 (Abel's theorem). *Suppose $\sum c_n$ converges. Let*

$$f(x) = \sum_{n=0}^{\infty} c_n x^n \quad (-1 < x < 1).$$

$$\text{Then } \lim_{x \rightarrow 1} f(x) = \sum_{n=0}^{\infty} c_n.$$

Proof. Let

$$s_n = c_0 + \cdots + c_n, \quad s_{-1} = 0.$$

Then

$$\begin{aligned} \sum_{n=0}^m c_n x^n &= \sum_{n=0}^m (s_n - s_{n-1}) x^n \\ &= (1-x) \sum_{n=0}^{m-1} s_n x^n + s_m x^m. \end{aligned}$$

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 for all $n \geq N$,

$$|s - s_n| < \frac{\varepsilon}{2}.$$

If $x > 1 - \delta$, for some suitably chosen $\delta > 0$, we have

$$\begin{aligned} |f(x) - s| &= \left| (1-x) \sum_{n=0}^{\infty} (s_n - s) x^n \right| \\ &\leq (1-x) \sum_{n=0}^N |s_n - s| |x|^n + \frac{\varepsilon}{2} \\ &\leq \varepsilon. \end{aligned}$$

Hence $\lim_{x \rightarrow 1} f(x) = s$, as desired. □

We now require a theorem concerning an inversion in the order of summation.

Proposition 13.9. *Given a double sequence (a_{ij}) , $i = 1, 2, 3, \dots$, $j = 1, 2, 3, \dots$, suppose that*

$$\sum_{j=1}^{\infty} |a_{ij}| = b_i \quad (i = 1, 2, 3, \dots)$$

and $\sum b_i$ converges. Then

$$\sum_{i=1}^{\infty} \sum_{j=1}^{\infty} a_{ij} = \sum_{j=1}^{\infty} \sum_{i=1}^{\infty} a_{ij}. \quad (13.1)$$

Theorem 13.10 (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

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x - a)^n \quad (|x - a| < R - |a|). \quad (13.2)$$

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 13.11. Suppose $\sum a_n x^n$ and $\sum b_n x^n$ converge in $(-R, R)$. Let E be the set of all $x \in (-R, R)$ 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 $(-R, R)$, then $a_n = b_n$ for $n = 0, 1, 2, \dots$. Hence (20) holds for all $x \in (-R, R)$.

Exponential and Logarithmic Functions

Definition 13.12 (Exponential function). For $z \in \mathbb{C}$, define

$$\exp(z) := \sum_{n=0}^{\infty} \frac{z^n}{n!}. \quad (13.3)$$

Proposition 13.13. $\exp(z)$ converges for every $z \in \mathbb{C}$.

Proof. Ratio test. □

Proposition 13.14. For $z, w \in \mathbb{C}$,

$$\exp(z + w) = \exp(z) \exp(w).$$

Corollary 13.15. For $z \in \mathbb{C}$,

$$\exp(z) \exp(-z) = 1.$$

Proof. Take $z = z, w = -z$ in the previous result. □

Proposition 13.16. \exp is strictly increasing in \mathbb{R} .

Proposition 13.17. For $z \in \mathbb{C}$,

$$\exp'(z) = \exp(z)$$

Further,

$$\exp'(z) = \lim_{h \rightarrow 0} \frac{\exp(z + h) - \exp(z)}{h} = \lim_{h \rightarrow 0} \frac{\exp(z + h) - 1}{h} \exp(z).$$

Let $\exp(1) = e$. So $\exp(n) = \exp(1 + \cdots + 1) = \exp(1) \cdots \exp(1) = e^n$. This holds for any $n \in \mathbb{Q}$.

Trigonometric Functions

Definition 13.18. For $z \in \mathbb{C}$, define

$$\begin{aligned}\cos z &:= \frac{1}{2} (\exp(iz) + \exp(-iz)) \\ \sin z &:= \frac{1}{2i} (\exp(iz) - \exp(-iz))\end{aligned}\tag{13.4}$$

By Eq. (13.3), we obtain the power series

$$\begin{aligned}\cos z &= \sum_{n=0}^{\infty} (-1)^n \frac{z^{2n}}{(2n)!} \\ \sin z &= \sum_{n=0}^{\infty} (-1)^n \frac{z^{2n+1}}{(2n+1)!}\end{aligned}$$

We now state some properties of trigonometric functions:

- (Euler's identity) $\exp(iz) = \cos z + i \sin z$.
- For $x \in \mathbb{R}$, $\cos' x = -\sin x$ and $\sin' x = \cos x$.

Proposition 13.19.

- (i) \exp is periodic, with period $2\pi i$.
- (ii) C and S are periodic, with period 2π .
- (iii) If $0 < t < 2\pi$, then $\exp(it) \neq 1$.
- (iv) If $z \in \mathbb{C}$, $|z| = 1$, there exists a unique $t \in [0, 2\pi)$ such that $\exp(it) = z$.

§13.2 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 13.20 (Fundamental Theorem of Algebra). *Suppose a_0, \dots, a_n are complex numbers, $n \geq 1$, $a_n \neq 0$,*

$$P(z) = \sum_{k=0}^n a_k z^k.$$

Then $P(z) = 0$ for some complex number z .

Proof. WLOG assume $a_n = 1$. Let

$$\mu = \inf |P(z)| \quad (z \in \mathbb{C}).$$

If $|z| = R$, then

$$|P(z)| \geq R^n (1 - |a_{n-1}|R^{-1} - \dots - |a_0|R^{-n}).$$

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

$$Q(z) = 1 + b_k z^k + \dots + b_n z^n \quad (b_k \neq 0).$$

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

$$|Q(re^{i\theta})| \leq 1 - r^k (|b_k| - r|b_{k+1}| - \dots - r^{n-k}|b_n|).$$

For sufficiently small r , the expression in braces is positive; hence $|Q(re^{i\theta})| < 1$, a contradiction.

Thus $\mu = 0$, that is, $P(z_0) = 0$. □

§13.3 Fourier Series

Definition 13.21. A *trigonometric polynomial* is a finite sum of the form

$$f(x) = a_0 + \sum_{n=1}^{\infty} (a_n \cos nx + b_n \sin nx) \quad (x \in \mathbb{R})$$

where $a_0, \dots, a_N, b_1, \dots, b_N \in \mathbb{C}$.

We can write the above in the form

$$f(x) = \sum_{n=-N}^N c_n e^{inx}.$$

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 13.22 (Fourier coefficients). If f is an integrable function on $[-\pi, \pi]$, the numbers c_m defined by

$$c_m = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) e^{inx} dx$$

for all integers m are called the *Fourier coefficients* of f .

Definition 13.23 (Fourier series). The series

$$\sum_{n=-\infty}^{\infty} c_n e^{inx}$$

formed with the Fourier coefficients is called the *Fourier series* of f .

Definition 13.24. Let (ϕ_n) be a sequence of complex functions on $[a, b]$ such that

$$\int_a^b \phi_n(x) \overline{\phi_m(x)} dx = 0 \quad (n \neq m).$$

Then (ϕ_n) is said to be an *orthogonal system of functions* on $[a, b]$. If in addition

$$\int_a^b |\phi_n(x)|^2 dx = 1 \quad (n = 1, 2, \dots)$$

then (ϕ_n) is said to be *orthonormal*.

Example 13.25. The functions $\phi_n = \frac{1}{\sqrt{2\pi}} e^{inx}$ form an orthonormal system on $[-\pi, \pi]$.

If (ϕ_n) is orthonormal on $[a, b]$ and if

$$c_n = \int_a^b f(t) \overline{\phi_n(t)} dt \quad (n = 1, 2, 3, \dots)$$

we call c_n the n -th Fourier coefficient of f relative to (ϕ_n) . We write

$$f(x) \sim \sum_{n=1}^{\infty} c_n \phi_n(x)$$

and call this series the Fourier series of f (relative to (ϕ_n)).

Remark. The symbol \sim used above implies nothing about the convergence of the series; it merely says that the coefficients are given by (66).

Proposition 13.26. Let (ϕ_n) be orthonormal 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 assume

$$t_n(x) = \sum_{k=1}^n \gamma_k \phi_k(x).$$

Then

$$\int_a^b |f - s_n|^2 dx \leq \int_a^b |f - t_n|^2 dx$$

and 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 .

Proposition 13.27 (Bessel inequality). If (ϕ_n) is orthonormal on $[a, b]$, and if

$$f(x) \sim \sum_{n=1}^{\infty} c_n \phi_n(x)$$

then

$$\sum_{n=1}^{\infty} |c_n|^2 \leq \int_a^b |f(x)|^2 dx. \quad (13.5)$$

In particular, $c_n \rightarrow 0$.

Theorem 13.28 (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

$$\lim_{N \rightarrow \infty} \|f - s_N(f)\|_2^2 = \lim_{N \rightarrow \infty} \frac{1}{2\pi} \int_{-\pi}^{\pi} |f(x) - s_N(f; x)|^2 dx = 0.$$

Also

$$\langle f, g \rangle = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) \overline{g(x)} dx = \sum_{n=-\infty}^{\infty} c_n \overline{\gamma_n}$$

and

$$\|f\|_2^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} |f(x)|^2 dx = \sum_{n=-\infty}^{\infty} |c_n|^2.$$

§13.4 Gamma Function

Definition 13.29 (Gamma function). For $0 < x < \infty$, the *Gamma function* is defined as

$$\Gamma(x) := \int_0^\infty t^{x-1} e^{-t} dt. \quad (13.6)$$

The integral converges for these x . (When $x < 1$, both 0 and ∞ have to be looked at.)

Lemma 13.30.

(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.

(ii) Since $\Gamma(1) = 1$, (1) implies (2) by induction.

(iii)

□

In fact, these three properties characterise Γ completely.

Lemma 13.31 (Characteristic properties 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 13.32 (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 13.33.

$$B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}.$$

Proof. Let $f(x) = \frac{\Gamma(x+y)}{\Gamma(y)}B(x, y)$. We want to prove that $f(x) = \Gamma(x)$, using Lemma 13.31.

(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 $\log f(x)$ is convex.

Therefore $f(x) = \Gamma(x)$ 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

$$\begin{aligned} \Gamma(x)\Gamma(y) &= 4 \int_0^{\frac{\pi}{2}} \left[\int_0^{+\infty} r^{2x-1} \cos^{2x-1} \theta \cdot r^{2y-1} \sin^{2y-1} \theta \cdot e^{-r^2} \cdot r dr \right] d\theta \\ &= \underbrace{2 \int_0^{\frac{\pi}{2}} \cos^{2x-1} \theta \sin^{2y-1} \theta d\theta}_{B(x,y)} \cdot \underbrace{2 \int_0^{+\infty} r^{2(x+y)-1} e^{-r^2} dr}_{\Gamma(x+y)} \end{aligned}$$

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:

$$I = \int_{-\infty}^{+\infty} e^{-x^2} dx \quad (13.7)$$

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 13.34.

$$\Gamma(x) = \frac{2^{x-1}}{\sqrt{\pi}} \Gamma\left(\frac{x}{2}\right) \Gamma\left(\frac{x+1}{2}\right).$$

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 13.35 (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*

$$\lim_{x \rightarrow \infty} \frac{\Gamma(x+1)}{(x/e)^x \sqrt{2\pi x}} = 1. \quad (13.8)$$

Proof. □**Lemma 13.36.**

$$B(p, 1-p) = \Gamma(p) \Gamma(1-p) = \frac{\pi}{\sin p\pi}.$$

Proof. □

Exercises

V

General Topology

The study of topology simultaneously simplifies and generalises the theory of metric spaces. By discarding the metric, and focusing solely on the more basic and fundamental notion of an open set, many arguments and proofs are simplified. And many constructions (such as the important concept of a quotient space) cannot be carried out in the setting of metric spaces: they need the more general framework of topological spaces.

14 Topological Spaces and Continuous Functions

§14.1 Topological Spaces

Definition 14.1 (Topological space). A *topology* on a set X is a collection \mathcal{T} of subsets of X satisfying:

- (i) $X, \emptyset \in \mathcal{T}$;
- (ii) if $U_i \in \mathcal{T}$ for all $i \in I$, then $\bigcup_{i \in I} U_i \in \mathcal{T}$; (closed under arbitrary unions)
- (iii) if $U_1, \dots, U_n \in \mathcal{T}$, then $\bigcap_{i=1}^n U_i \in \mathcal{T}$. (closed under finite intersections)

A set X for which a topology \mathcal{T} has been specified is called a *topological space*, denoted by (X, \mathcal{T}) . $U \subset X$ is called an *open set* of X if $U \in \mathcal{T}$.

Notation. When \mathcal{T} is understood, we talk about the topological space X .

Example 14.2. Let X be any non-empty set.

- The *discrete topology* on X is the set of all subsets of X ; that is, $\mathcal{T} = \mathcal{P}(X)$.
- The *indiscrete topology* (or *trivial topology*) on X is $\mathcal{T} = \{X, \emptyset\}$.
- The *co-finite topology* on X consists of the empty set together with every subset U of X such that $X \setminus U$ is finite.
- Let \mathcal{T}_c be the collection of all subsets $U \subset X$ such that U^c either is countable or is all of X . Then \mathcal{T}_c is a topology on X .

Definition 14.3. Suppose \mathcal{T} and \mathcal{T}' are two topologies on a given set X . We say that

- (i) \mathcal{T} is *finer* than \mathcal{T}' if $\mathcal{T} \supset \mathcal{T}'$;
- (ii) \mathcal{T} is *coarser* than \mathcal{T}' if $\mathcal{T} \subset \mathcal{T}'$;
- (iii) \mathcal{T} is *comparable* with \mathcal{T}' if either $\mathcal{T} \supset \mathcal{T}'$ or $\mathcal{T} \subset \mathcal{T}'$.

Remark. The indiscrete topology is the coarsest topology possible, while the discrete topology is the finest topology possible.

§14.2 Basis for a Topology

In linear algebra, every vector space is generated by a basis. In topology, we have a similar notion, as it is usually hard to define a topology by specifying all the open sets.

Definition 14.4 (Basis). A **basis** for a topology on X is a collection \mathcal{B} of subsets of X (called *basis elements*) if

- (i) for all $x \in X$, there exists $B \in \mathcal{B}$ such that $x \in B$;
- (ii) for all $B_1, B_2 \in \mathcal{B}$ and $x \in B_1 \cap B_2$, there exists $B_3 \in \mathcal{B}$ such that $x \in B_3 \subset B_1 \cap B_2$.

We define the topology \mathcal{T} generated by basis \mathcal{B} as

$$U \in \mathcal{T} \iff \forall x \in U, \exists B \in \mathcal{B}, x \in B \subset U. \quad (14.1)$$

We now check that the collection \mathcal{T} generated by the basis \mathcal{B} is, in fact, a topology on X .

- (i) \emptyset satisfies the defining condition of openness vacuously, so $\emptyset \in \mathcal{T}$. $X \in \mathcal{T}$ follows from (i) of Definition 14.4.

- (ii) Consider a collection $\{U_i \mid i \in I\}$ of elements of \mathcal{T} . We want to show that $U = \bigcup_{i \in I} U_i \in \mathcal{T}$.

Given $x \in U$, there exists $i \in I$ such that $x \in U_i$. Since $U_i \in \mathcal{T}$, there exists $B \in \mathcal{B}$ such that $x \in B \subset U_i$. Thus $x \in B \subset U$, so $U \in \mathcal{T}$.

- (iii) Take two elements $U_1, U_2 \in \mathcal{T}$, we want to show that $U_1 \cap U_2 \in \mathcal{T}$.

Given $x \in U_1 \cap U_2$, choose $B_1 \in \mathcal{B}$ such that $x \in B_1 \subset U_1$; choose $B_2 \in \mathcal{B}$ such that $x \in B_2 \subset U_2$. Then $x \in B_1 \cap B_2$.

Since \mathcal{B} is a basis, by (ii) of Definition 14.4, there exists $B_3 \in \mathcal{B}$ such that $x \in B_3 \subset B_1 \cap B_2$. Thus $U_1 \cap U_2 \in \mathcal{T}$.

Finally, we show by induction that any finite intersection $U_1 \cap \cdots \cap U_n \in \mathcal{T}$. This is trivial for $n = 1$; suppose it true for $n - 1$ and prove it for n . Now

$$(U_1 \cap \cdots \cap U_n) = (U_1 \cap \cdots \cap U_{n-1}) \cap U_n.$$

By hypothesis, $U_1 \cap \cdots \cap U_{n-1} \in \mathcal{T}$. Thus by the result just proved, the intersection of $U_1 \cap \cdots \cap U_{n-1}$ and U_n also belongs to \mathcal{T} .

Another way of describing the topology generated by a basis is given in the following result:

Lemma 14.5. Let \mathcal{T} be the topology on X generated by basis \mathcal{B} . Then \mathcal{T} equals the collection of all unions of elements of \mathcal{B} .

Proof. Let $\mathcal{B} = \{B_i \mid i \in I\}$.

- If $B_i \in \mathcal{B}$, see that

$$\forall x \in B, x \in B \subset B \implies B \in \mathcal{T}.$$

Since \mathcal{T} is a topology, the arbitrary unions of B_i 's must be in \mathcal{T} .

- Conversely, given $U \in \mathcal{T}$, for each $x \in U$, there exists $B_x \in \mathcal{B}$ such that $x \in B_x \subset U$. Then $U = \bigcup_{x \in U} B_x$, so U is a union of elements of \mathcal{B} .

□

Remark. The above result states that every $U \in \mathcal{T}$ can be expressed as a union of basis elements.

We have described in two different ways how to go from a basis to the topology it generates. Sometimes we need to go in the reverse direction, from a topology to a basis generating it. Here is one useful way of obtaining a basis for a given topology.

Lemma 14.6. *Let (X, \mathcal{T}) be a topological space. Suppose that \mathcal{C} is a collection of open sets of X , such that*

$$\forall U \in \mathcal{T}, \quad \forall x \in U, \quad \exists C \in \mathcal{C}, \quad x \in C \subset U.$$

Then \mathcal{C} is a basis for \mathcal{T} .

Proof. We first show that \mathcal{C} is a basis.

- (i) For all $x \in X$, since $X \in \mathcal{T}$, by hypothesis, there exists $C \in \mathcal{C}$ such that $x \in C \subset X$.
- (ii) Let $x \in C_1 \cap C_2$, where $C_1, C_2 \in \mathcal{C} \subset \mathcal{T}$. Thus $C_1, C_2 \in \mathcal{T}$, so $C_1 \cap C_2 \in \mathcal{T}$. Hence by hypothesis, there exists $C_3 \in \mathcal{C}$ such that $x \in C_3 \subset C_1 \cap C_2$.

Let \mathcal{T}' be the topology generated by \mathcal{C} . We will show that $\mathcal{T} = \mathcal{T}'$.

- Let $U \in \mathcal{T}, x \in U$. By hypothesis, there exists $C \in \mathcal{C}$ such that $x \in C \subset U$. By definition, $U \in \mathcal{T}'$. Hence $\mathcal{T} \subset \mathcal{T}'$.
- Conversely, let $W \in \mathcal{T}'$. By Lemma 14.5, W is a union of elements of \mathcal{C} . Since each element of \mathcal{C} is an element of \mathcal{T} (and thus open), and a union of open sets is open, so $W \in \mathcal{T}$. Hence $\mathcal{T}' \subset \mathcal{T}$.

□

When topologies are given by bases, the next result is a criterion to determine whether one topology is finer than another.

Lemma 14.7. *Let \mathcal{B} and \mathcal{B}' be bases for the topologies \mathcal{T} and \mathcal{T}' respectively on X . Then the following are equivalent:*

- (i) \mathcal{T}' is finer than \mathcal{T} .
- (ii) For all $x \in X$, and for all $B \in \mathcal{B}$ such that $x \in B$, there exists $B' \in \mathcal{B}'$ such that $x \in B' \subset B$.

Proof.

(ii) \implies (i) Let $U \in \mathcal{T}$. To show that $\mathcal{T} \subset \mathcal{T}'$, we want to show that $U \in \mathcal{T}'$.

Let $x \in U$. Since \mathcal{B} generates \mathcal{T} , there exists $B \in \mathcal{B}$ such that $x \in B \subset U$. By (ii), there exists $B' \in \mathcal{B}'$ such that $x \in B' \subset B$. Then $x \in B' \subset U$, so $U \in \mathcal{T}'$, by definition.

(i) \implies (ii) We are given $x \in X$ and $B \in \mathcal{B}$, with $x \in B$.

Now $B \in \mathcal{T}$ by definition, and $\mathcal{T} \subset \mathcal{T}'$ by (i); therefore, $B \in \mathcal{T}'$. Since \mathcal{T}' is generated by \mathcal{B}' , there exists $B' \in \mathcal{B}'$ such that $x \in B' \subset B$. \square

We now define three topologies on the real line \mathbb{R} .

Definition 14.8.

- (i) Let \mathcal{B} be the collection of all open intervals in \mathbb{R} . The topology generated by \mathcal{B} is called the **standard topology** on \mathbb{R} .

Whenever we consider \mathbb{R} , we shall suppose it is given this topology unless stated otherwise.

- (ii) Let \mathcal{B}' be the collection of all half-open intervals of the form $[a, b)$. The topology generated by \mathcal{B}' is called the **lower limit topology** on \mathbb{R} .

When \mathbb{R} is given the lower limit topology, we denote it by \mathbb{R}_ℓ .

- (iii) Let $K = \{\frac{1}{n} \mid n \in \mathbb{Z}^+\}$, and let \mathcal{B}'' be the collection of all open intervals (a, b) , along with all sets of the form $K \setminus (a, b)$. The topology generated by \mathcal{B}'' is called the **K -topology** on \mathbb{R} .

When \mathbb{R} is given this topology, we denote it by \mathbb{R}_K .

It is easy to see that all three of these collections are bases; in each case, the intersection of two basis elements is either another basis element or is empty. The relation between these topologies is the following:

Lemma 14.9. *The topologies of \mathbb{R}_ℓ and \mathbb{R}_K are strictly finer than the standard topology on \mathbb{R} , but are not comparable with one another.*

Definition 14.10 (Subbasis). A **subbasis** \mathcal{S} for a topology on X is a collection of subsets of X whose union equals X .

The **topology \mathcal{T} generated by the subbasis \mathcal{S}** is defined as the collection of all unions of finite intersections of elements of \mathcal{S} :

$$U \in \mathcal{T} \iff U = \text{union of finite intersections in } \mathcal{S}.$$

We now check that \mathcal{T} is a topology. Consider the collection

$$\mathcal{B} = \{\text{all finite intersections of elements of } \mathcal{S}\}.$$

It suffices to show that \mathcal{B} is a basis, for then by Lemma 14.5, the collection \mathcal{T} of all unions of elements of \mathcal{B} is a topology.

- (i) Given $x \in X$, it belongs to an element of \mathcal{S} and hence to an element of \mathcal{B} .

- (ii) Let

$$B_1 = S_1 \cap \cdots \cap S_m, \quad B_2 = S'_1 \cap \cdots \cap S'_n$$

be two elements of \mathcal{B} . Their intersection

$$B_1 \cap B_2 = (S_1 \cap \cdots \cap S_m) \cap (S'_1 \cap \cdots \cap S'_n)$$

is also a finite intersection of elements of \mathcal{S} , so it belongs to \mathcal{B} .

§14.3 Examples of Topologies

Order Topology

Definition 14.11 (Order topology). Let $(X, <)$, $|X| > 1$. Let \mathcal{B} be the collection of all sets of the following types:

- (i) All open intervals (a, b) in X .
- (ii) All intervals of the form $[a_0, b)$, where a_0 is the smallest element (if any) of X .
- (iii) All intervals of the form $(a, b_0]$, where b_0 is the largest element (if any) of X .

The topology generated by \mathcal{B} is called the *order topology*.

We need to check that \mathcal{B} is a basis of X .

- (i) Every $x \in X$ lies in some element of \mathcal{B} : the smallest element (if any) lies in all sets of type (ii), the largest element (if any) lies in all sets of type (iii), and every other element lies in a set of type (i).
- (ii) The intersection of any two sets of the preceding types is a set of one of these types, or is empty. Several cases need to be checked; we leave it to you.

For instance, let $x \in (a, b) \cap (c, d)$. Let $p = \max\{a, c\}$, $q = \min\{b, d\}$. Then $x \in (p, q) \subset (a, b) \cap (c, d)$, where $(p, q) \in \mathcal{B}$.

Example 14.12. • The standard topology on \mathbb{R} is just the order topology derived from the usual order on \mathbb{R} .

Definition 14.13. Let $(X, <)$, $a \in X$. Then the following subsets of X are *rays* determined by a :

$$\begin{aligned} (a, +\infty) &= \{x \in X \mid x > a\}, \\ [a, +\infty) &= \{x \in X \mid x \geq a\}, \\ (-\infty, a) &= \{x \in X \mid x < a\}, \\ (-\infty, a] &= \{x \in X \mid x \leq a\}. \end{aligned}$$

$(a, +\infty)$ and $(-\infty, a)$ are called *open rays*, since they are open; for instance, $(a, +\infty) = \bigcup_{x>a} (a, x)$. Similarly, $[a, +\infty)$ and $(-\infty, a]$ are *closed rays*.

Lemma 14.14. The collection of open rays form a subbasis for the order topology.

Proof. Let \mathcal{T} be the order topology on X , let \mathcal{T}' be the topology generated by the subbasis of open rays. We will show that $\mathcal{T} = \mathcal{T}'$.

- Because the open rays are open in the order topology, the topology they generate is contained in the order topology. Hence $\mathcal{T}' \subset \mathcal{T}$.
- On the other hand, every basis element for the order topology equals a finite intersection of open rays; the interval (a, b) equals the intersection of $(-\infty, b)$ and $(a, +\infty)$, while $[a_0, b)$ and $(a, b_0]$, if they exist, are themselves open rays. Hence the topology generated by the open rays contains the order topology, so $\mathcal{T} \subset \mathcal{T}'$.

□

Product Topology

Definition 14.15. Let (X, \mathcal{T}_X) and (Y, \mathcal{T}_Y) be topological spaces. The *product topology* on $X \times Y$ is the topology $\mathcal{T}_{X \times Y}$ with basis

$$\mathcal{B} = \{U \times V \mid U \in \mathcal{T}_X, V \in \mathcal{T}_Y\}.$$

We first check that \mathcal{B} is a basis.

- (i) $X \times Y$ is a basis element, so every element of $X \times Y$ is contained in $X \times Y$.
- (ii) Let $U_1 \times V_1, U_2 \times V_2 \in \mathcal{B}$. Then their intersection is

$$(U_1 \times V_1) \cap (U_2 \times V_2) = (U_1 \cap U_2) \times (V_1 \cap V_2).$$

Since $U_1 \cap U_2 \in \mathcal{T}_X$, $V_1 \cap V_2 \in \mathcal{T}_Y$, we have that $(U_1 \cap U_2) \times (V_1 \cap V_2) \in \mathcal{B}$.

Subspace Topology

Definition 14.16 (Subspace). Let (X, \mathcal{T}) be a topological space. If $Y \subset X$, the collection

$$\mathcal{T}_Y := \{Y \cap U \mid U \in \mathcal{T}\}$$

is a topology on Y , called the *subspace topology*. With this topology, Y is called a *subspace* of X ; its open sets consist of all intersections of open sets of X with Y .

We check that \mathcal{T}_Y is a topology.

Lemma 14.17. If \mathcal{B} is a basis for the topology of X , then

$$\mathcal{B}_Y = \{B \cap Y \mid B \in \mathcal{B}\}$$

is a basis for the subspace topology on Y .

Lemma 14.18. Let Y be a subspace of X . If U is open in Y , and Y is open in X , then U is open in X .

Proposition 14.19. *If A is a subspace of X , and B is a subspace of Y , then the product topology on $A \times B$ is the same as the topology $A \times B$ inherits as a subspace of $X \times Y$.*

§14.4 Closed Sets and Limit Points

Let X be a topological space.

Note that if U is an open set containing x , we often say that U is a *neighbourhood* of x .

Closed Sets

Definition 14.20 (Closed set). $A \subset X$ is *closed* if its complement A^c is open.

The collection of closed subsets of a space X has properties similar to those satisfied by the collection of open subsets of X :

Lemma 14.21. *Let X be a topological space.*

- (i) \emptyset and X are closed.
- (ii) Arbitrary intersections of closed sets are closed.
- (iii) Finite unions of closed sets are closed.

Proof.

(i) \emptyset and X are closed because they are the complements of the open sets X and \emptyset , respectively.

(ii) Suppose $\{A_i \mid i \in I\}$ is a collection of closed sets. By de Morgan's laws,

$$\left(\bigcap_{i \in I} A_i \right)^c = \bigcup_{i \in I} A_i^c.$$

Since A_i^c 's are open, the RHS is open since it is an arbitrary union of open sets. Hence $\bigcap A_i$ is closed.

(iii) Suppose A_i is closed for $i = 1, \dots, n$. Then

$$\left(\bigcup_{i=1}^n A_i \right)^c = \bigcap_{i=1}^n A_i^c.$$

The RHS is a finite intersection of open sets and is thus open. Hence $\bigcup A_i$ is closed.

□

Remark. Note that \emptyset and X are both open and closed. This explains the statement “a door is not a set”: a door must be either open or closed, and cannot be both, while a set can be open, or closed, or both, or neither!

If Y is a subspace of X , we say A is closed in Y if $A \subset Y$ and A is closed in the subspace topology of Y (that is, if $Y \setminus A$ is open in Y). We have the following result:

Proposition 14.22. *Let Y be a subspace of X . Then A is closed in Y if and only if it equals the intersection of a closed set of X with Y .*

Proof.

\Leftarrow Assume that $A = C \cap Y$, where C is closed in X . Then $X \setminus C$ is open in X , so that $(X \setminus C) \cap Y$ is open in Y , by definition of the subspace topology. But $(X \setminus C) \cap Y = Y \setminus A$. Hence $Y \setminus A$ is open in Y , so that A is closed in Y .

\Rightarrow Suppose A is closed in Y . Then $Y \setminus A$ is open in Y , so that by definition it equals the intersection of an open set U of X with Y . The set $X \setminus U$ is closed in X , and $A = Y \cap (X \setminus U)$, so that A equals the intersection of a closed set of X with Y , as desired. \square

Proposition 14.23. *Let Y be a subspace of X . If A is closed in Y , and Y is closed in X , then A is closed in X .*

Closure and Interior

Definition 14.24. The *interior* of $A \subset X$ is the union of all open sets contained in A , denoted by $\text{Int } A$.

The *closure* of A is the intersection of all closed sets contained in A , denoted by \overline{A} .

Proposition 14.25. *Let Y be a subspace of X ; let $A \subset Y$, let \overline{A} denote the closure of A in X . Then the closure of A in Y equals $\overline{A} \cap Y$.*

Limit Points

Definition 14.26. Suppose $A \subset X$. $x \in X$ is a limit point of A if every neighbourhood of x intersects A in some point other than x itself.

A' denotes the set of all limit points of A .

Proposition 14.27. *Let $A \subset X$. Then $\overline{A} = A \cup A'$.*

Corollary 14.28. *$A \subset X$ is closed if and only if it contains all its limit points.*

Hausdorff Spaces

Definition 14.29 (Hausdorff space). A **Hausdorff space** is a topological space X such that for all distinct $x_1, x_2 \in X$, there exist neighbourhoods U_1 and U_2 of x_1 and x_2 respectively that are disjoint.

Proposition 14.30. *Every finite point set in a Hausdorff space X is closed.*

The condition that finite point sets be closed is in fact weaker than the Hausdorff condition. For example, \mathbb{R} in the finite complement topology is not a Hausdorff space, but it is a space in which finite point sets are closed. The condition that finite point sets be closed has been given a name of its own: it is called the *T1 axiom*.

Proposition 14.31. *Let X be a space satisfying the T1 axiom; let $A \subset X$. Then x is a limit point of A if and only if every neighborhood of x contains infinitely many points of A .*

Proposition 14.32. *If X is a Hausdorff space, then a sequence of points of X converges to at most one point of X .*

Proposition 14.33. *Every simply ordered set is a Hausdorff space in the order topology. The product of two Hausdorff spaces is a Hausdorff space. A subspace of a Hausdorff space is a Hausdorff space.*

Bibliography

- [Apo57] T. M. Apostol. *Mathematical Analysis*. Addison-Wesley, 1957.
- [Art11] M. Artin. *Algebra*. Pearson Education, 2011.
- [Axl24] S. Axler. *Linear Algebra Done Right, 4th edition*. 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.
- [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.

Index

- analytic function, 288
- annihilator, 122
- balls, 153
 - closed ball, 153
 - open ball, 153
 - punctured ball, 153
- basis, 84
- beta functions, 298
- boundary, 158
- boundary point, 158
- boundedness, 154
- Cauchy sequence, 188
- closed set, 157
- closure, 158
- compact, 163
 - open cover, 163
- connectedness, 176
- continuity, 214
 - uniform continuity, 221
- convergence, 180
- coset, 56, 118
 - left coset, 56
 - right coset, 56
- Dedekind cut, 133
- dense, 158
- dimension, 87
- direct sum, 78
- dual basis, 121
- dual map, 122
- equivalence relation, 26
 - equivalence class, 26
 - partition, 26
 - quotient set, 27
- extended real number system, 142
- finite-dimensional, 81
- Fourier coefficients, 295
- Fourier series, 295
- function, 31
 - bijectivity, 33
 - image, 31
 - injectivity, 33
 - invertibility, 36
 - pre-image, 31
 - restriction, 31
 - surjectivity, 33
- Gamma function, 298
- group, 48
- homomorphism, 62
- image, 63, 96
- index, 57
- induced set, 160
- infimum, 127
- injectivity, 96
- interior, 158
- invertibility, 108
- isomorphism, 62, 109
- kernel, 63, 96
- limit of function, 211
- limit point, 160
- linear combination, 80
- linear functional, 121
- linear independence, 81
- linear map, 93
- matrix, 101
 - identity matrix, 113
 - transpose, 106
- matrix of linear map, 101
- matrix of vector, 111
- metric space, 152
- neighbourhood, 155
- open set, 155
- order, 127
- orthogonal system of functions, 295
- perfect set, 173
- pointwise convergence, 265

- power series, 287
- product of vector spaces, 116
- quotient map, 119
- quotient space, 118
- rank, 107
 - column rank, 106
 - column space, 106
 - row rank, 106
 - row space, 106
- relation, 25
 - binary relation, 25
 - partial order, 26
 - total order, 26
 - well order, 26
- Riemann–Stieltjes integrability, 246
- set, 20
 - Cartesian product, 22
 - complement, 22
 - disjoint, 22
 - element, 20
 - empty set, 20
 - intersection, 22
 - interval, 21
 - ordered pair, 22
 - power set, 21
 - set difference, 22
 - subset, 21
 - union, 22
- span, 80
- subgroup, 52
- subsequence, 186
- supremum, 127
- surjectivity, 97
- uniform convergence, 267
- vector space, 72
 - complex vector space, 72
 - real vector space, 72
 - subspace, 76