# Analysis, Topology, & Manifolds

Mario L. Gutierrez Abed

# Contents

# Preface

As of now, this book is merely a compilation of course notes that I took while I was an undergraduate student at CUNY Hunter College. As such, large portions of the text is copied nearly verbatim from the literature cited on the references (something I intend to correct later on). I am currently at a stage in my career in which I have to focus all my energy on doing research rather than writing books; writing this one was merely a "fun exercise." Nevertheless, I have rewritten many of these notes with my own flavor and there are some original proofs here and there as well as exercises from homework problems that I had to solve as part of my own coursework. My main goal is to provide undergraduate students with a meaty summary of some of the major topics that (I think) every undergraduate student should have before pursuing graduate studies. Of course, a textbook of such nature must necessarily sacrifice depth of content for scope; indeed there are plenty of textbooks out there dedicated to just one of these subjects alone! The main advantage of having a book like this one is to have a wide-scope reference available whenever you need to brush up your knowledge, and it is also a resource that may serve you well as a companion guide to your main textbooks. I sincerely hope that you enjoy the readings and learn something (or a lot!) from them.

# Chapter 1

# Real Analysis

REAL ANALYSIS is an exciting challenge of wrestling with big ideas related to the development of the mathematical machinery of real numbers and real-valued functions of real variables. You may have heard that this subject is "calculus on steroids;" this is indeed a fitting "definition" in my opinion (seriously, make sure you have some Tylenol[1] within reach before you start reading this; I kid you not). That being said however, once you get past the headaches you will find yourself immersed in a rich world of elegant mathematics at which point you will go sprinting to your registrar's office and declare yourself a math major if you haven't already. Getting through this course will not be an easy ride, but hey, nothing worth fighting for comes easy, right? Real analysis is one of your first true expositions to advanced mathematics, and as such you are bound to find numerous counterintuitive surprises. Having seen mainly graphical, numerical, or intuitive arguments throughout all your baby-math courses, you will now need to learn what constitutes a rigorous mathematical proof and how to write one. You will need to be convinced of the need for a more rigorous study of functions. The necessity of precise definitions and an axiomatic approach must be carefully motivated. There needs to be significant reward for the difficult work of firming up the logical structure of limits. Specifically, real analysis should not be just an elaborate reworking of standard introductory calculus. You will be exposed to the tantalizing complexities of the real line, to the subtleties of different flavors of convergence, and to the intellectual delights hidden in the paradoxes of the infinite. Get ready to have your mind blown over and over again!

---

[1] This is not an endorsement!

## 1.1  Real Number System

Our goal in this first section is to provide a quick review of a handful of important ideas from advanced calculus (and to encourage a bit of practice on these fundamentals). We will make no attempt to be thorough. Our purpose is to set the stage for later generalizations and to collect together in one place some of the notation that should already be more or less familiar.

**Proposition 1.** *Let $p$ be a prime number. Then there are no integers $m, n \in \mathbb{Z}$ such that*

$$\left(\frac{m}{n}\right)^2 = p. \tag{1.1}$$

*Proof.* Suppose to the contrary that such a number does in fact exist. We can assume that this fraction $m/n$ is in simplest terms, i.e. $\gcd(m, n) = 1$, since otherwise we could just divide by the common factor and get an equivalent fraction. Then the equation (1.1) can be rewritten as $m^2 = n^2 p$, which means that $p \mid m^2$ (this notation means that $p$ divides $m^2$). Thus, it follows that $p \mid m$. In particular, there exists an integer $k$ such that $m = kp$, from which it follows that

$$m^2 = k^2 p^2 = n^2 p \implies k^2 p = n^2.$$

But this would mean that $p \mid n^2$ and therefore that $p \mid n$. This is a contradiction, because $1 = \gcd(m, n) \geq p > 1$. Thus, we conclude that no rational function $x^2 = p$ exists. $\qquad\square$

Now we make the following observation: Let $p$ be a prime number, and set

$$A = \{r \in \mathbb{Q}^+ \mid r^2 < p\} \qquad \text{and} \qquad B = \{r \in \mathbb{Q}^+ \mid r^2 > p\}.$$

Then for each $r \in A$, there exists $s \in A$ such that $r < s$. Similarly, for each $r \in B$ there exists $s \in B$ such that $s < r$. Let us proceed to define this number $s$:

- If $r \in A$, then

$$r^2 < p \implies r < \sqrt{p} = r + (\sqrt{p} - r)$$

$$= r + (\sqrt{p} - r) \cdot \frac{(\sqrt{p} + r)}{(\sqrt{p} + r)}$$

$$= r + \frac{p - r^2}{\sqrt{p} + r} > \underbrace{r + \frac{p - r^2}{p + r}}_{\text{we call this number } s} > r.$$

- Similarly, if $r \in B$, we have

$$r^2 > p \implies r > \sqrt{p} = r - (r - \sqrt{p})$$
$$= r - (r - \sqrt{p}) \cdot \frac{(r + \sqrt{p})}{(r + \sqrt{p})}$$
$$= r - \frac{r^2 - p}{r + \sqrt{p}} < \underbrace{r - \frac{r^2 - p}{r + p}}_{\text{we call this number } s} < r.$$

Thus we have that $s \in \mathbb{Q}^+$. If $r \in A$, then $r^2 - p < 0$, implying that $r < s$. On the other hand, if $r \in B$, then $r^2 - p > 0$, implying that $r > s$.

**Remark**: The observation above suggests that any element in $B \subset \mathbb{Q}$ is an upper bound of $A$. In other words, if $s \in B$ and $r \in A$, then $r < s$. Furthermore, $A$ has no smallest upper bound (in $\mathbb{R}$): For any $s \in B$, there is an $s_1 < s$, with $s_1 \in B$, such that $s_1$ is an upper bound of $A$. Similar reasoning shows that $B$ is bounded below by elements in $A$ with no largest lower bound. We will soon examine this observation more closely.

**Definition 1.** *Let S be a set. An **order** on S is a relation, denoted by $<$, with the following two properties:*

- *If $x, y \in S$, then one and only one of the statements below is true:*

$$x < y, \qquad x = y, \qquad y < x.$$

- *Let $x, y, z \in S$. Then if $x < y$ and $y < z$, it is always true that $x < z$.*

*An **ordered set** is a set S in which an order is defined (for example, $\mathbb{Q}$ is an ordered set if $r < s$ is defined to mean that $s - r$ is a positive rational number). There are sets in which certain elements are ordered, but some are not. These sets are said to be **partially ordered** (for a silly example take for instance the set $\{1, 3, 4, pizza\}$. Then we can put an order (either increasing or decreasing) on the numbers $1, 3, 4$, but obviously that delicious pizza is left out of any possible ordering!).*

**Definition 2.** *Suppose S is an ordered set, and $E \subset S$. If there exists a $\beta \in S$ such that $x \leq \beta$ for every $x \in E$, then we say that E is bounded above, and call $\beta$ an **upper bound** of E. **Lower bounds** are defined in the same way (with $\geq$ in place of $\leq$).*

**Definition 3.** *A partially ordered set P is said to be a **directed set** if every pair of elements of P has an upper bound.*

**Definition 4.** *Suppose S is an ordered set, $E \subset S$, and E is bounded above. Moreover, suppose there exists an $\alpha \in S$ with the following properties:*

- *$\alpha$ is an upper bound of E.*

- *If $\gamma < \alpha$, then $\gamma$ is not an upper bound of E.*

*Then $\alpha$ is called the **least upper bound** or **supremum** of E, denoted $\alpha = \sup E$. The **greatest lower bound**, or **infimum**, of a set E which is bounded below is defined in the same manner. The statement $\alpha = \inf E$ means that $\alpha$ is a lower bound of E and that no $\beta$ with $\beta > \alpha$ is a lower bound of E.*

**Definition 5.** *If both $\sup(a, b)$ and $\inf(a, b)$ exist for every pair $(a, b)$ of a partially ordered set P, then P is said to be a **lattice**.*

**Example 1.** *a) Consider the sets A and B described above as subsets of the ordered set $\mathbb{Q}$. The set A is bounded above. In fact, the upper bounds of A are exactly the members of B. Since B contains no smallest member, A has no least upper bound in $\mathbb{Q}$. Similarly, B is bounded below: the set of all lower bounds of B consists of A and of all $r \in \mathbb{Q}$ with $r \leq 0$. Since A has no largest member, B has no greatest lower bound in $\mathbb{Q}$.*

*b) If $\alpha = \sup E$ exists, then $\alpha$ may or may not be a member of E. For instance, let $E_1$ be the set of all $r \in \mathbb{Q}$ with $r < 0$. Let $E_2$ be the set of all $r \in \mathbb{Q}$ with $r \leq 0$. Then $\sup E_1 = \sup E_2 = 0$ with $0 \notin E_1$ and $0 \in E_2$.*

*c) Let E consist of all numbers $1/n$, for $n \in \mathbb{Z}^+$. Then $\sup E = 1$, which is in E, and $\inf E = 0$, which is not in E.*

An ordered set $S$ is said to have the *least upper bound property* if the following is true: If $E \subset S$ is not empty, and $E$ is bounded above, then $\sup E$ exists in $S$. (Observe that $\mathbb{Q}$ does not have the least upper bound property.) We now show that every set $S$ with the least upper bound property also has the greatest lower bound property:

**Theorem 1.** *Suppose S is an ordered set with the least upper bound property and let $B \subset S$ be nonempty and bounded below. In addition, let L be the set of all lower bounds of B. Then $\alpha = \sup L$ exists in S and $\alpha = \inf B$. In particular, $\inf B$ exists in S.*

*Proof.* Since $B$ is bounded below, $L$ is nonempty. Since $L$ consists of exactly those $y \in S$ which satisfy the inequality $y \leq x$ for all $x \in B$, we see that every $x \in B$ is an upper bound of $L$. Then $L$ is bounded above. Our hypothesis about $S$ implies therefore that $L$ has a supremum in $S$; call it $\alpha$. If $\gamma < \alpha$, then $\gamma$ is not an upper bound of $L$. In particular, there is some $\beta \in L$ such that $\gamma < \beta$, implying that $\gamma$ is a lower bound of $B$. Thus $\alpha \leq x$ for all $x \in B$. It follows that $\alpha \in L$. If $\alpha < \lambda$, then $\lambda \notin L$, since $\alpha$ is an upper bound of $L$. Thus we have shown that $\alpha \in L$ but $\lambda \notin L$ if $\alpha < \lambda$. In other words, $\alpha$ is a lower bound of $B$ but $\lambda$ is not if $\lambda > \alpha$. This implies that $\alpha = \inf B$. $\qquad\qquad\square$

**Theorem 2 (Existence Theorem).** *There exists an ordered field $\mathbb{R}$ which has the least upper bound property. Moreover, $\mathbb{R}$ contains $\mathbb{Q}$ as a subfield.*

*Proof.* See [Rudin, 1964, Chapter 1, Appendix]. $\qquad\qquad\square$

We now derive some important properties of the field $\mathbb{R}$:

**Axiom 1 (Axiom of Completeness).** *Every nonempty set of real numbers that is bounded above has a least upper bound.*

**Theorem 3.** *Let $x, y \in \mathbb{R}$. Then,*

    *a) if $x > 0$, there is a positive integer n such that $nx > y$.*

    *b) if $x < y$, there exists a $p \in \mathbb{Q}$ such that $x < p < y$.*

*(Part a) is usually referred to as the Archimedian property of $\mathbb{R}$. Part b) may be stated by saying that $\mathbb{Q}$ is dense in $\mathbb{R}$: Between any two real numbers there is a rational one.)*

*Proof of a).* Set $A = \{nx \mid n \in \mathbb{N}\}$. Now let's assume that a) is false, so that $y$ is an upper bound of $A$, and put $\alpha = \sup A$. We have that $x > 0$, which implies that $\alpha - x < \alpha$, and this

in turn means that $\alpha - x$ is not an upper bound of $A$. Hence $\alpha - x < mx$ for some positive integer $m$. But then $\alpha < (m+1)x \in A$, which contradicts the statement that $\alpha = \sup A$. Therefore $A$ is not bounded above. $\qquad\square$

*Proof of b).* Since $x < y$, we have $y - x > 0$. From part a), we conclude that there is an integer $n > 0$ such that $n(y - x) > 1$. Observe that, for some integer $m$, we have $m - 1 \le nx < m$. Observe also that $m \le 1 + nx < ny$. Thus, since $nx < m$, we have $nx < m < ny$. In particular, $x < m/n < y$. This proves b), with $p = m/n$. $\qquad\square$

Now we are ready to prove the existence of $n^{\text{th}}$ roots of positive reals, but before doing so let us recall the binomial theorem:

$$(x + y)^n = \sum_{k=0}^{n} \binom{n}{k} x^k y^{n-k} \qquad \text{where} \qquad \binom{n}{k} = \frac{n!}{k!(n-k)!}.$$

Let us illustrate this with two simple examples:

$$a) \quad (x+y)^2 = y^2 \binom{2}{0} + xy \binom{2}{1} + x^2 \binom{2}{2}$$
$$= y^2 \frac{2!}{0!\,2!} + xy \frac{2!}{1!\,1!} + x^2 \frac{2!}{2!\,0!}$$
$$= y^2 + 2xy + x^2.$$

$$b) \quad (x+y)^3 = y^3 \binom{3}{0} + xy^2 \binom{3}{1} + x^2 y \binom{3}{2} + x^3 \binom{3}{3}$$
$$= y^3 \frac{3!}{0!\,3!} + xy^2 \frac{3!}{1!\,2!} + x^2 y \frac{3!}{2!\,1!} + x^3 \frac{3!}{3!\,0!}$$
$$= y^3 + 3xy^2 + 3x^2 y + x^3.$$

**Theorem 4.** *For every real $x > 0$ and every integer $n > 0$, there is one and only one real $y$ such that $y^n = x$.*

*Proof.* One proof is already given in [Rudin, 1964, p. 10]. I offer here an alternate proof, which relies on the binomial theorem. Let $E$ be defined as in Rudin's; that is, we let $E$ be

the set consisting of all positive real numbers $t$ such that $t^n < x$, i.e., $E = \{t \in \mathbb{R}^+ \mid t^n < x\}$, with $x \in \mathbb{R}$. Check that $E$ is not empty and that it is bounded (see Rudin's argument if you wish). Then set $y = \sup E$. We will show that $y^n = x$ by proving that $\|y^n - x\| < \varepsilon$ for any $\varepsilon > 0$. This will imply that $\|y^n - x\| = 0$ or $y^n = x$.

Let $h > 0$. If $h < y$, then $y - h$ is a positive number that is not an upper bound of $E$. In particular, there is some $t \in E$ such that $y - h < t < y$ and therefore $(y - h)^n < t^n < x$, from which follows that $y - h \in E$. Observe now that $(y + h)^n > x$, for if $(y + h)^n \leq x$, then $(y + h/2)^n < (y + h)^n \leq x$, implying that $y + h/2 \in E$, which contradicts the fact that $y$ is an upper bound of $E$. It follows that $(y - h)^n < x < (y + h)^n$ and naturally, $(y - h)^n < y^n < (y + h)^n$.

Geometrically, the distance from $y^n$ to $x$, given by $\|y^n - x\|$, is less than $(y + h)^n - (y - h)^n$. This can be proven analytically without much difficulty. Thus, we have

$$\|y^n - x\| < (y + h)^n - (y - h)^n = 2 \sum_{i=0}^{\lfloor \frac{n-1}{2} \rfloor} \binom{n}{2i + 1} y^{n-2i-1} h^{2i+1} \qquad (\clubsuit)$$

$$< 2h \sum_{i=0}^{\lfloor \frac{n-1}{2} \rfloor} \binom{n}{2i + 1} y^{n-2i-1}.$$

(Here we use $\lfloor - \rfloor$ to denote the usual *floor function*) The expression on the right hand side of $(\clubsuit)$ was derived from expanding $(y + h)^n - (y - h)^n$ with the help of the binomial theorem. The last inequality was derived from the assumption that $h$ may be selected to be less than 1.

Now notice that

$$B = 2 \sum_{i=0}^{\lfloor \frac{n-1}{2} \rfloor} \binom{n}{2i + 1} y^{n-2i-1}.$$

is just a number and $hB < \varepsilon$ for any $\varepsilon$ given a sufficiently small $h$. Thus, $\|y^n - x\| < \varepsilon$, as desired. $\qquad \square$

**Corollary 1.** *If $a$ and $b$ are positive real numbers and $n$ is a positive integer, then $(ab)^{1/n} = a^{1/n} b^{1/n}$.*

*Proof.* Let $\alpha = a^{1/n}$ and $\beta = b^{1/n}$, so that

$$ab = \alpha^n \beta^n = (\alpha \beta)^n,$$

since multiplication is commutative. The uniqueness assertion of Theorem 4 shows that

$$(ab)^{1/n} = \alpha\beta = \alpha^{1/n}\beta^{1/n}. \qquad \square$$

One approach to describe the elements of $\mathbb{R}$ is by using decimals; the following propositions give some insight. First, however, we must discuss the mathematical concept of a *sequence*. You most likely have a fairly accurate intuition of what a sequence is; we can think of it as a "list" of elements that may or may not have a particular order assigned to such elements. However, in order for sequences to be useful to us, we need a more precise mathematical definition:

**Definition 6.** *A **sequence** $\{x_n\}$ (also usually denoted as $(x_n)$) of points in a set X is a mapping $\mathbb{N} \to X$ by $n \mapsto x_n$.*

**Remark:** For the most part, we are mainly interested in sequences. However, note that in some situations you may also encounter mappings that generalize the concept of sequences: a **net** $\{x_\alpha\}_{\alpha \in \mathfrak{D}}$ (or just $\{x_\alpha\}$, or $(x_\alpha)$) is a mapping $\mathfrak{D} \to X$ by $\alpha \mapsto x_\alpha$, where $\mathfrak{D}$ is a directed set that may be "uncountable" (this is in contrast to sequences, where the domain $\mathbb{N}$ of such mappings is "countable."[2]

**Proposition 2.** *Fix an integer $p \geq 2$ and let $\{a_n\}$ be any sequence of integers satisfying $0 \leq a_n \leq p-1$ for all n. Then, $\sum_{n=1}^{\infty} a_n/p^n$ converges to a number in $[0,1]$.*

*Proof.* Since $a_n \geq 0$, the partial sums $\sum_{n=1}^{N} a_n/p^n$ are nonnegative and increase with $N$. Thus, to show that the series converges to some number in $[0,1]$, we just need to show that 1 is an upper bound for the sequence of partial sums:

$$\sum_{n=1}^{N} \frac{a_n}{p^n} \leq \sum_{n=1}^{N} \frac{p-1}{p^n} \leq (p-1)\sum_{n=1}^{\infty} \frac{1}{p^n} = 1. \qquad \square$$

Consequently, each $x$ in $[0,1]$ can be so represented:

---

[2] These concepts of "countable" and "uncountable" sets are discussed in the next subsection. You will then see why it's important to consider both *sequences* and *nets*.

**Proposition 3.** *Let $p \in \mathbb{Z}$ such that $p \geq 2$, and let $x \in [0, 1]$. Then there is a sequence of integers $\{a_n\}$ with $0 \leq a_n \leq p - 1$ for all $n$ such that $x = \sum_{n=1}^{\infty} a_n / p^n$.*

*Proof.* The case $x = 0$ is trivial (this is a word that mathematicians like to throw around all the time, although in cases like this it is certainly justified!) so let us suppose that $0 < x \leq 1$. We will then construct $\{a_n\}$ by induction. Start by choosing $a_1$ to be the largest integer satisfying $a_1 / p < x$. Since $x > 0$, it follows that $a_1 \geq 0$, and since $x \leq 1$, we have $a_1 < p$. Now because $a_1$ is an integer, this means that $a_1 \leq p - 1$. Also, since $a_1$ is largest, we must have

$$\frac{a_1}{p} < x \leq \frac{a_1 + 1}{p}.$$

Next, choose $a_2$ to be the largest integer satisfying

$$\frac{a_1}{p} + \frac{a_2}{p^2} < x.$$

Check that $0 \leq a_2 \leq p - 1$ and that

$$\frac{a_1}{p} + \frac{a_2}{p^2} < x \leq \frac{a_1}{p} + \frac{a_2 + 1}{p^2}.$$

Thus, by induction we get a sequence of integers $\{a_n\}$ with $0 \leq a_n \leq p - 1$ such that

$$\frac{a_1}{p} + \cdots + \frac{a_n}{p^n} < x \leq \frac{a_1}{p} + \cdots + \frac{a_n + 1}{p^n}.$$

Now it follows that $x = \sum_{n=1}^{\infty} a_n / p^n$, and we are done.                    $\square$

The series $\sum_{n=1}^{\infty} a_n / p^n$ is called a *base $p$* (or *$p$-adic*) decimal expansion for $x$. It is sometimes written in the shorter form

$$x = 0.a_1 a_2 a_3 \ldots \text{(base } p).$$

It does not have to be unique (even for ordinary base-10 decimals: $0.5 = 0.4999\ldots$). One problem is that our construction is designed to produce nonterminating decimal expansions. In the particular case where

$$x = \frac{a_1}{p} + \cdots + \frac{a_n + 1}{p^n} = \frac{q}{p^n} \qquad \text{for some integer } 0 < q \leq p^n,$$

the construction will give us a repeating string of $p - 1$'s in the decimal expansion for $x$ since

$$\frac{1}{p^n} = \sum_{k=n+1}^{\infty} \frac{p-1}{p^k}.$$

That is, any such $x$ has two distinct base $p$ decimal expansions:

$$x = \frac{a_1}{p} + \cdots + \frac{a_n + 1}{p^n} = \frac{a_1}{p} + \cdots + \frac{a_n}{p^n} + \sum_{k=n+1}^{\infty} \frac{p-1}{p^k}.$$

Notice that if $y \in \mathbb{R}$, for any $n \in \mathbb{N}$ we have $y \in [n, n+1]$. In particular, there is some $x \in [0,1]$ such that $y = n + x$. By the work done above, this means that any real number $y$ is an infinite sum of rational numbers.

## 1.2   Finite, Countable, and Uncountable Sets

If $A$ and $B$ are sets and there exists a one-to-one mapping of $A$ onto $B$ (that is, a map from $A$ to $B$ that is bijective), then we say that $A$ and $B$ can be put in $1 - 1$ correspondence, or that $A$ and $B$ have the same cardinal number (intuitively, the same "number of elements." We will discuss this at greater length below). It is not hard to show that cardinality yields an *equivalence relation*: if $A$ and $B$ have the same cardinal number, we say that $A$ and $B$ are equivalent, and we write $A \sim B$. This relation clearly has the following properties:

- It is reflexive: $A \sim A$.

- It is symmetric: If $A \sim B$, then $B \sim A$.

- It is transitive: For some other set $C$, if $A \sim B$ and $B \sim C$, then $A \sim C$.

Any relation with these three properties is called an ***equivalence relation***.

**Definition 7.** *For any positive integer n, let* $\mathbb{N}_{(n)} \subset \mathbb{N}$ *be the set whose elements are the integers* $1, 2, \ldots, n$ *(*$\mathbb{N}$ *of course represents the set of all natural numbers; it is the extension of* $\mathbb{N}_{(n)}$ *to infinity). Then, for any set A we say:*

    *i)* A is **finite** if $A \sim \mathbb{N}_{(n)}$ for some *n* (the empty set is also considered to be finite).

    *ii)* A is **infinite** if A is not finite. (Duh!)

    *iii)* A is **countable** if $A \sim \mathbb{N}$.

    *iv)* A is **uncountable** if A is neither finite nor countable.

    *v)* A is **at most countable** if A is finite or countable.

*Countable sets are sometimes called **enumerable**, or **denumerable**.*

    For two finite sets *A* and *B*, we evidently have $A \sim B$ if and only if *A* and *B* contain the same number of elements. For infinite sets however, the idea of "having the same number of elements" becomes quite vague, whereas the notion of $1 - 1$ correspondence retains its clarity.

**Example 2.** *a) The set of all integers $\mathbb{Z}$ is countable. To see this we can define a function $f : \mathbb{Z} \to \mathbb{N}$ such that*

$$f(n) = \begin{cases} 2n & \text{if } n \geq 1, \\ -2n + 1 & \text{if } n \leq 0. \end{cases}$$

*This function sets up the $1 - 1$ correspondence*

$$
\begin{array}{ccccccccc}
\mathbb{Z}: & \ldots & -3, & -2, & \ldots, & 2, & 3, & \ldots \\
 & & \updownarrow & \updownarrow & & \updownarrow & \updownarrow & \\
\mathbb{N}: & \ldots & 7, & 5, & \ldots, & 4, & 6, & \ldots
\end{array}
$$

*That is, the negative integers in $\mathbb{Z}$ are mapped to the odd numbers on $\mathbb{N}$ while the positive integers in $\mathbb{Z}$ are mapped to the even numbers in $\mathbb{N}$.*

    *Note that usually there are multiple bijective maps capable of establishing a $1 - 1$ correspondence between two sets. For instance, we could have instead used a map from $\mathbb{N}$ to $\mathbb{Z}$, say $g \colon \mathbb{N} \to \mathbb{Z}$ such that*

$$g(n) = \begin{cases} \frac{n}{2} & \text{if } n \text{ is even,} \\ -\frac{n-1}{2} & \text{if } n \text{ is odd.} \end{cases}$$

*This function sets up the $1-1$ correspondence*

$$\mathbb{N}: \quad \ldots \quad 1, \quad 2, \quad 3, \quad 4, \quad \ldots$$
$$\updownarrow \quad \updownarrow \quad \updownarrow \quad \updownarrow$$
$$\mathbb{Z}: \quad \ldots \quad 0, \quad 1, \quad -1, \quad 2, \quad \ldots$$

   *That both $f$ and $g$ are bijective is easy to check. Notice that $\mathbb{Z}$ is equivalent to a proper subset of itself!*



You, completely mind-blown.

   *This is typical of infinite sets whereas it is impossible for finite sets.*

   *b) The set of all cartesian products on $\mathbb{N}$ is equivalent to $\mathbb{N}$ itself, i.e., $\mathbb{N} \times \mathbb{N} \sim \mathbb{N}$. A quick proof is supplied by the fundamental theorem of arithmetic: Each positive integer $k \in \mathbb{N}$ can be uniquely written as $k = 2^{m-1}(2n-1)$ for some $m, n \in \mathbb{N}$. Define $f: \mathbb{N} \times \mathbb{N} \to \mathbb{N}$ by $f(m, n) = 2^{m-1}(2n-1)$. It is obvious that this $f$ is bijective.*

   *c) The set of all real numbers $\mathbb{R}$ is equivalent to the interval $(-\pi/2, \pi/2)$. To see this, define $f: \mathbb{R} \to (-\pi/2, \pi/2)$ by $f(x) = \tan^{-1}(x)$. Recall from calculus that $f$ is a strictly increasing (hence one-to-one) function from $\mathbb{R}$ to $(-\pi/2, \pi/2)$, and it is also clearly onto. (As a matter of fact, a generalization of this exercise is that $\mathbb{R}$ is in fact equivalent to any interval of real numbers $(a, b)$. Mind-blowing indeed, isn't it??)*

**Theorem 5.** *Every infinite subset of a countable set $A$ is countable.*

*Proof.* Suppose $E \subset A$ and $E$ is infinite. Arrange the elements $x$ of $A$ in a sequence $\{x_n\}$ of distinct elements. Then construct a sequence $\{n_k\}$ as follows: Let $n$ be the smallest positive integer such that $x_{n_1} \in E$. Having chosen $n_1, \ldots, n_{k-1}$ for $k = 2, 3, 4, \ldots$, let $n_k$ be the smallest integer greater than $n_{k-1}$ such that $x_{n_k} \in E$. Now putting $f(k) = x_{n_k}$ for $k = 1, 2, 3, \ldots$, we obtain a $1 - 1$ correspondence between $E$ and $\mathbb{N}$. $\qquad\square$

If $A$ and $B$ are sets, then it is customary to denote the set of all points in $A$ that are not in $B$ by $A \smallsetminus B$; that is, $A \smallsetminus B = \{x \in A \mid x \notin B\}$. That being said, we have the following theorem:

**Theorem 6.** *Every infinite set has a countable subset.*

*Proof.* Let $A$ be an infinite set. Then $A \neq \varnothing$, because $\varnothing$ is considered to be finite. Let $x_1 \in A$ be any element of $A$. Then $A \smallsetminus \{x_1\} \neq \varnothing$ (otherwise $A = \{x_1\}$ and $A$ is finite). Pick $x_2 \in A \smallsetminus \{x_1\}$ to be any element of $A \smallsetminus \{x_1\}$. Having chosen $x_1, \ldots, x_{n-1}$, observe that $A \smallsetminus \{x_1, \ldots, x_{n-1}\} \neq \varnothing$ (otherwise $A = \{x_1, \ldots, x_{n-1}\}$, making $A$ finite). Hence we are free to select $x_n \in A \smallsetminus \{x_1, \ldots, x_{n-1}\}$. Let $E = \{x_n\} \subset A$. Then $E$ is countable. $\qquad\square$

Theorem 6 above shows that a countable infinity is the smallest type of infinity. That is, no uncountable set can be a subset of a countable set, while every infinite set has a countable subset. To motivate our next several results, we now present a second proof that $\mathbb{N} \times \mathbb{N}$ is equivalent to $\mathbb{N}$.

**Theorem 7.** $\mathbb{N} \times \mathbb{N}$ *is equivalent to* $\mathbb{N}$.

*Proof.* Arrange $\mathbb{N} \times \mathbb{N}$ in a matrix as in Figure 1.1. The arrows have been added to show how we are going to enumerate $\mathbb{N} \times \mathbb{N}$. We will count the pairs in the order indicated by the arrows: $(1, 1)$, $(2, 1)$, $(1, 2)$, $(3, 1)$, $(2, 2)$, and so on, accounting for each upward slanting diagonal in succession.

Notice that all of the pairs along a given diagonal have the same sum. The entries of $(1, 1)$ add to 2, the entries of both $(2, 1)$ and $(1, 2)$ add to 3, each pair of entries on the next diagonal add to 4, and so on. Moreover, for any given $n$, there are exactly $n$ pairs whose

$$
\begin{array}{cccc}
(1,1) & (1,2) & (1,3) & (1,4) \quad \cdots \\
(2,1) & (2,2) & (2,3) & \cdots \\
(3,1) & (3,2) & \vdots \\
(4,1) & \vdots
\end{array}
$$

Figure 1.1: $\mathbb{N} \times \mathbb{N}$ arranged as a matrix.

entries sum to $n + 1$. In other words, there are exactly $n$ pairs on the $n^{th}$ diagonal. These observations, allow us to construct an explicit formula for this correspondence between $\mathbb{N} \times \mathbb{N}$ and $\mathbb{N}$: Let $f \colon \mathbb{N} \times \mathbb{N} \to \mathbb{N}$ be defined by

$$
f(m,n) = \frac{(m + n - 2)(m + n - 1)}{2} + n.
$$

This is an invertible (i.e., bijective) map from $\mathbb{N} \times \mathbb{N}$ to $\mathbb{N}$; hence $\mathbb{N} \times \mathbb{N} \sim \mathbb{N}$, and this proves our theorem.                                                                                 □

The above theorem gives us a ton of new information. We can see this materialize in the following theorem:

**Theorem 8.** *The countable union of countable sets is countable. That is, if $A_i$ is countable for $i = 1, 2, 3, \ldots$, then $\cup_{i=1}^{\infty} A_i$ is countable.*

*Proof.* Since each $A_i$ is countable, we can arrange their elements collectively in a matrix:

$$
\begin{array}{llllll}
A_1 : & a_{11} & a_{12} & a_{13} & . & . \\
A_2 : & a_{21} & a_{22} & a_{23} & . & . \\
A_3 : & a_{31} & a_{32} & a_{33} & . & . \\
& . & . & . & . & . & . \\
& . & . & . & . & . & .
\end{array}
$$

So $\cup_{i=1}^{\infty} A_i$ is the range of some invertible map on $\mathbb{N} \times \mathbb{N}$ (just as the one constructed on the Theorem 7 above). That is, $\cup_{i=1}^{\infty} A_i$ is equivalent to $\mathbb{N} \times \mathbb{N}$ and hence to $\mathbb{N}$.                                    □

Note that proof of the above theorem can be used to show that, given any two countable sets $A$ and $B$, the set $A \times B$ is also countable.

**Corollary 2.** $\mathbb{Q}$ *is countable.*

Recall that between any two real numbers there is a rational number (i.e., $\mathbb{Q}$ is dense in $\mathbb{R}$). This means, in fact, that between any two real numbers, there are infinitely many rational numbers (since $\mathbb{R}$ is infinite and we know that $(a,b) \sim \mathbb{R}$ for all $a, b \in \mathbb{R}$). Surprisingly, $\mathbb{N}$ is as large as $\mathbb{Q}$ even though $\mathbb{N} \subset \mathbb{Q}$ and there are infinitely many rationals between any two rational numbers.

So far we have shown that $\mathbb{N}, \mathbb{Z}, \mathbb{Q}$ are all countable. Now we show the shocking result that $\mathbb{R}$ is not a countable set.

**Theorem 9.** $\mathbb{R}$ *is uncountable.*

*Proof.* To prove that $\mathbb{R}$ is uncountable, it is enough to show that some subset in $\mathbb{R}$ is uncountable  (since no countable set can have an uncountable subset). Therefore, we can use the subset $(0,1)$ and prove that it is uncountable. To accomplish this, we will show that any countable subset of $(0,1)$ is proper.

Given any sequence $\{a_n\}$ in $(0,1)$, we construct an element $x$ in $(0,1)$ with $x \neq a_n$ for any $n$. We begin by listing the decimal expansions of the $a_n$; for example:

$$a_1 = 0.\boxed{3}\,1572\ldots$$
$$a_2 = 0.0\,\boxed{4}\,268\ldots$$
$$a_3 = 0.91\,\boxed{5}\,36\ldots$$
$$a_4 = 0.759\,\boxed{9}\,9\ldots$$
$$\ldots = \ldots\ldots\ldots\ldots$$

(If any $a_n$ has two representations, just use the infinite one.) Now let $x = 0.533353\ldots$, where the $n^{\text{th}}$ digit in the expansion for $x$ is taken to be 3, unless $a_n$ happens to have 3 as its $n^{\text{th}}$ digit, in which case we replace it with 5 (this is why we "boxed" the $n^{\text{th}}$ digit in the expansion of $a_n$ above. Note that the choices of 3 and 5 are more or less arbitrary, in truth we just want to avoid the troublesome digits 0 and 9 but any other digits would do).

Using this procedure, the decimal representation of $x$ is unique because it does not end in all 0's or all 9's, and $x \neq a_n$ for any $n$ because the decimal expansions for $x$ and $a_n$ differ in the $n^{\text{th}}$ place. Thus we have shown that $\{a_n\}$ is a proper subset of $(0, 1)$, and hence $(0, 1)$ is uncountable, which in turn implies that $\mathbb{R}$ is uncountable. $\qquad\square$

The proof that we just produced is known as ***Cantor's diagonalization method***. It gives insight into the differences between countable and uncountable sets.

**Corollary 3.** *The set of all irrationals $\mathbb{R} \setminus \mathbb{Q}$ (or simply $\mathbb{I}$), is uncountable.*

*Proof.* We know that $\mathbb{R} = \mathbb{Q} \cup \mathbb{I}$. We also know that the union of countable sets must be countable. Since $\mathbb{Q}$ is countable and $\mathbb{R}$ is uncountable, it follows that $\mathbb{I}$ must be uncountable, as desired. $\qquad\square$

**Theorem 10 (Cantor's Theorem).** *Let $A$ be a set and $\mathcal{P}(A)$ denote its power set (i.e., the collection of all subsets of $A$ including $A$ itself). Then no map $F\colon A \to \mathcal{P}(A)$ can be onto.*

*Proof.* Let us assume that $F\colon A \to \mathcal{P}(A)$ is an onto function. Then consider $S_F = \{x \in A : x \notin F(x)\} \in \mathcal{P}(A)$. Since $F$ is assumed to be onto, there must exists an element $y \in A$ such that $F(y) = S_F$. We claim that $S_F \neq F(y)$ for any $y \in A$. Indeed, if $S_F = F(y)$, then we are faced with the following alternatives:

$$y \in F(y) = S_F \implies y \notin F(y) \qquad \text{or} \qquad y \notin F(y) = S_F \implies y \in F(y),$$

and both lead to contradictions! $\qquad\square$

**Theorem 11 (Bernstein's Theorem).** *Let $A$ and $B$ be nonempty sets. If there exist one-to-one maps $f\colon A \to B$ and $g\colon B \to A$, then there is a map $h\colon A \to B$ that is both one-to-one and onto. Informally, this implies that if two cardinalities are both less than or equal to each other, then they are equal.*

*Proof.* One proof can be found in [Carothers, 2000]. We will present here an alternate proof. We will call an element $b \in B$ *lonely* if there is no element $a \in A$ such that $f(a) = b$. We also say that an element $b_1$ of $B$ is a *descendent* of an element $b_0 \in B$ if there is a natural

number $n$ (possibly zero) such that $b_1 = (f \circ g)^n (b_0)$. We now define the function $h \colon A \to B$ as follows:

$$h(a) = \begin{cases} g^{-1}(a) & \text{if } f(a) \text{ is the descendant of a lonely point,} \\ f(a) & \text{otherwise.} \end{cases}$$

Note that if $f(a)$ is the descendent of a lonely point, then $f(a) = f \circ g(b)$ for some $b$; since $g$ is injective, the element $g^{-1}(a)$ is well defined. Thus our function $h$ is well defined. We claim that it is a bijection from $A$ to $B$.

We first prove that $h$ is surjective. Indeed, if $b \in B$ is the descendent of a lonely point, then $h(g(b)) = b$; and if $b$ is not the descendent of a lonely point, then $b$ is not lonely, so there is some $a \in A$ such that $f(a) = b$; by our definition, then, $h(a) = b$. Thus it is surjective.

Next, we prove that $h$ is injective. We first note that for any $a \in A$, the point $h(a)$ is a descendent of a lonely point if and only if $f(a)$ is a descendent of a lonely point. Now suppose that we have two elements $a_1, a_2 \in A$ such that $h(a_1) = h(a_2)$. We consider two cases:

- If $f(a_1)$ is the descendent of a lonely point, so is $f(a_2)$. Then we have

$$g^{-1}(a_1) = h(a_1) = h(a_2) = g^{-1}(a_2).$$

  Since $g$ is a well defined function, it follows that $a_1 = a_2$.

- On the other hand, if $f(a_1)$ is not a descendent of a lonely point, then neither is $f(a_2)$. It follows that

$$f(a_1) = h(a_1) = h(a_2) = f(a_2).$$

  Since $f$ is injective, we have $a_1 = a_2$, which in turn implies that $h$ is injective as well.

Thus we have shown that $h$ is both surjective and injective, and hence bijective. $\qquad\square$

To appreciate how incredible Bernstein's result truly is, consider the following example:

**Example 3.** *Let $\mathbb{R}^\infty$ be the set of all real-valued sequences. That is, if $x \in \mathbb{R}^\infty$, then*

$$x = (x_1, x_2, \ldots, x_n, \ldots), \qquad \text{where each } x_i \in \mathbb{R}.$$

*Then $\mathbb{R}^\infty \sim (0,1)$. To show this, first observe that $\mathbb{R}^\infty \sim (0,1)^\infty$. This observation is justified by defining a map $f \colon \mathbb{R}^\infty \to (0,1)^\infty$ by*

$$f(x_1, x_2, \dots) = \left( \frac{\tan^{-1}(x_1) + \frac{\pi}{2}}{\pi}, \frac{\tan^{-1}(x_2) + \frac{\pi}{2}}{\pi}, \dots \right).$$

*Thus, it is enough to show that $(0,1) \sim (0,1)^\infty$ (Note that $(0,1)^\infty$ is the set of all sequences $\{x_n\}$ with $x_n \in (0,1)$ ). To do this, observe that $f \colon (0,1) \to (0,1)^\infty$ given by $f(x) = (x,0,0,\dots)$ (the choice of zeroes is arbitrary, what is important is to fix the first element) is an injective map from $(0,1)$ into $(0,1)^\infty$. Thus,*

$$\mathrm{card}(0,1) \leq \mathrm{card}(0,1)^\infty$$

*To prove the other direction, let $x \in (0,1)^\infty$. Then $x = (x_1, \dots, x_n, \dots)$, where $x_n \in (0,1)$ for all $n \in \mathbb{N}$. Represent each $x_n$ by its unique finite decimal expansion*

$$x_n = 0.x_{n_1} x_{n_2} x_{n_3} \dots.$$

*In addition, let $p_n$ be the $n^{th}$ prime and define $g \colon (0,1)^\infty \to (0,1)$ by $g(x) = 0.y_1 y_2 y_3 \dots$, where*

$$y_k = \begin{cases} x_{n_i} & \text{if } k = p_n^i, \\ 0 & \text{otherwise.} \end{cases}$$

*Then $g$ is injective. In particular,*

$$\mathrm{card}(0,1)^\infty \leq \mathrm{card}(0,1).$$

*Thus it follows, by Bernstein's theorem, that $\mathbb{R}^\infty \sim (0,1)^\infty \sim (0,1)$.* ✿

**Theorem 12.** *The rational numbers $\mathbb{Q}$ have measure 0 (i.e. occupy no space) on the real number line.*[3]

*Proof.* Since $\mathbb{Q}$ is a countable set, we can list all of its elements in a sequence $\{x_n\}$. We will show that $\mathbb{Q}$ has measure 0 by proving that for any $\varepsilon > 0$, there is a collection of open intervals which cover $\mathbb{Q}$ and whose combined length is less than $\varepsilon$. To do this, for each $x_n \in \mathbb{Q}$, define $I_n$ by

$$I_n = \left( x_n - \frac{\varepsilon}{2^{n+1}}, x_n + \frac{\varepsilon}{2^{n+1}} \right).$$

---

[3] We will see more on measure theory on Section §1.9.

In other words, $I_n$ is just an interval of length $L(I_n) = \varepsilon/2^n$ centered at $x_n$. Clearly, it is true that $\mathbb{Q} \subset \cup_{n=1}^{\infty} I_n$. Thus we have

$$L(\cup_{n=1}^{\infty} I_n) \leq \sum_{n=1}^{\infty} L(I_n) = \sum_{n=1}^{\infty} \frac{\varepsilon}{2^n} = \varepsilon \sum_{n=1}^{\infty} \frac{1}{2^n} = \varepsilon. \qquad \square$$

Theorem 12 can be interpreted as saying that the likelihood of selecting a rational number at random in the set of real numbers is 0. To put it in more colorful terms, having selected one object, the chance that another randomly selected object can be described in terms of the first is 0.

**Definition 8.** *A number is said to be **algebraic** if there exist integers $a_0, a_1, \ldots, a_n \in \mathbb{Z}$ such that $a_0 + a_1 x + \cdots + a_n x^n = 0$.*

**Theorem 13.** *The set of all algebraic numbers is countable.*

*Proof.* Let $A_n$ be the set of all polynomials of degree $n$ with integer coefficients. The map $a_0 + a_1 x + \cdots + a_n x^n \mapsto (a_0, a_1, \ldots, a_n)$ shows that $A_n \sim \mathbb{Z}^{n+1}$, which implies that $A_n$ is countable. Now the set of all polynomials with integer coefficients can be written as the countably infinite union $A = \cup_{n=1}^{\infty} A_n$, which must therefore be countable. Thus, each polynomial in $A$ can be assigned a natural number that uniquely identifies it.

Let $k \in \mathbb{N}$ be the unique positive integer corresponding to $p(x) = a_0 + a_1 x + \cdots + a_n x^n$. Observe that this polynomial can have at most $n$ distinct complex roots. We can arrange these roots in lexicographic (i.e., dictionary) order from smallest to largest and associate $k.1$ with the smallest root of $p$, $k.01$ with the next smallest root of $p$, $k.001$ with the third smallest root, etc. Clearly, each algebraic number is thus paired with at least one rational number. This implies that algebraic numbers are countable. $\qquad \square$

Notice that all countable sets have measure 0 in $\mathbb{R}$ (or $\mathbb{C}$). Thus, the probability that a number is algebraic is 0, which implies that almost all numbers are transcendental! This is completely crazy talk! Think about how many transcendental numbers you can come up with off the top of your head....you can probably think of $\pi$, $e$, and a small handful others..., but according to our recent discovery, there are more of these weirdos than "regular" (i.e., algebraic) numbers! Mind-blown again!!

Your face right now.

## 1.3 Metrics and Norms

We will start this section with a discussion of *metrics* and *metric spaces*. This is an extremely important topic that will permeate much of the advanced mathematical machinery that awaits you in your future studies. The concept of a metric appears ubiquitously in mathematics, especially in analysis, topology, and geometry, as you will see on later chapters.

**Definition 9.** *Given a set $M$, a function $d\colon M \times M \to [0, \infty)$ satisfying the following properties is called a* **metric** *on $M$:*

  *i)* $d(x, y) = 0$ *if and only if $x = y$.*

  *ii)* (**Symmetry**) $d(x, y) = d(y, x)$ *for all pairs $x, y \in M$.*

  *iii)* (**Triangle Inequality**) $d(x, y) \leq d(x, z) + d(z, y)$ *for all $x, y, z \in M$.*

*The couple $(M, d)$, consisting of a set $M$ together with a metric $d$ defined on $M$, is called a* **metric space**.

**Example 4.** *a) Every set M admits at least one metric. For example, for any $x, y \in M$, the function d defined by*

$$d(x,y) = \begin{cases} 1 & \text{if } x \neq y, \\ 0 & \text{if } x = y, \end{cases}$$

*is a metric (check that the properties hold!). This mundane but always available metric is called the* **discrete metric** *on M. A set supplied with its discrete metric is called a* **discrete space**.

*b) An important example is the real line $\mathbb{R}$ together with its usual metric $d(x,y) = |x - y|$, where $|\cdot|$ denotes the usual absolute value. Any time we refer to $\mathbb{R}$ without explicitly naming a metric, the absolute value metric is always understood to be the one that we have in mind.*

*c) Any subset of a metric space is also a metric space in its own right. If d is a metric on M and $A \subseteq M$, then $d(x,y)$ is defined for any pair $x, y \in A$, i.e., $d|_{A \times A}(x,y) = d(x,y)$ whenever $x, y \in A$. Thus we will use the same letter d and simply refer to the metric space $(A, d)$ as opposed to $(A, d|_{A \times A})$. Of particular interest in this regard is that $\mathbb{N}, \mathbb{Z}, \mathbb{Q}$, and $\mathbb{R} \setminus \mathbb{Q}$ each come already supplied with a natural metric, namely the restriction of the usual absolute value metric of $\mathbb{R}$. In each case, we will refer to this restriction as the usual metric.* ✇

Alright, so we have thus far the discrete metric and the good ol' absolute value at our disposal. Now how can we enrich our arsenal of metric functions? Let us present a few options on the following example:

**Example 5.** *a) Suppose that d is a metric on M and $\varphi: M \to M$ is a bijection. Then $\Phi: M \times M \to \mathbb{R}$ defined by $\Phi(x,y) = d(\varphi(x), \varphi(y))$ is also a metric on M, as you can verify.*

*b) Suppose that $\psi: M \to \mathbb{R}$ is an injective function (not necessarily surjective). Then, $\Psi: M \times M \to \mathbb{R}$ defined by $\Psi(x,y) = |\psi(x) - \psi(y)|$ is a metric on M. This is an example of a* **pullback metric**; *the name comes from the fact that we are defining the metric on the domain space (M in this case) by "pulling back" the metric from the target space ($\mathbb{R}$ in this case). The concept of* **pullbacks** *is extremely important, and you will find it over and over again all across mathematics. Anyhow, to see that $\Psi$ is indeed a metric, note that it satisfies all the required properties:*

- $0 \leq \Psi(x,y) < \infty$ *for all pairs $x, y \in M$.*

- $\Psi(x,y) = 0$ *if and only if $|\psi(x) - \psi(y)| = 0$, which is true only if and only if $\psi(x) = \psi(y)$. But this happens only when $x = y$, since by assumption $\psi$ is injective.*

- $\Psi(x,y) = |\psi(x) - \psi(y)| = |\psi(y) - \psi(x)| = \Psi(y,x)$.

- *The triangle inequality is also satisfied: For all $x, y, z \in M$,*

$$\begin{aligned} \Psi(x,y) = |\psi(x) - \psi(y)| &= |\psi(x) - \psi(z) + \psi(z) - \psi(y)| \\ &\leq |\psi(x) - \psi(z)| + |\psi(z) - \psi(y)| \\ &= \Psi(x,z) + \Psi(z,y). \end{aligned}$$

*c) Define $d_1, d_2, d_3 \colon \mathbb{R}^2 \to [0, \infty)$ by*

$$d_1(x,y) = |\tan^{-1}(x) - \tan^{-1}(y)|, \quad d_2(x,y) = |x^3 - y^3|, \quad \text{and} \quad d_3(x,y) = |e^x - e^y|.$$

*Then, $d_1$, $d_2$, and $d_3$ are all metric functions on $\mathbb{R}^2$ (as you should check!).*

*d) Let $M = (0, \infty)$ and define $d_1, d_2, d_3 \colon M \times M \to [0, \infty)$ by*

$$d_1(x,y) = |\sqrt{x} - \sqrt{y}|, \quad d_2(x,y) = |\log x - \log y| = \left|\log \frac{x}{y}\right|, \quad \text{and} \quad d_3(x,y) = \left|\frac{1}{x} - \frac{1}{y}\right|.$$

*Then $d_1$, $d_2$, and $d_3$ are all metric functions on $M = (0, \infty)$ (as you should check!).*

*e) Note that a function can be a metric on one set and fail to be a metric on another. Take, for instance, the function $d(x,y) = |x^2 - y^2|$. Then $d$ defines a metric on $[0, \infty)$, but fails to be a metric on all of $\mathbb{R}$. We can easily see why, as it violates some properties of metric spaces (check the properties and you'll see!).* ✿

Now our arsenal of metrics has grown quite a bit, but why should we be satisfied when we can expand our collection of metrics even further? To do this, we first prove the following lemma:

**Lemma 1.** *Let $f \colon [0, \infty) \to [0, \infty)$ be any function with the following two properties:*

*a) $f(x) = 0 \iff x = 0$. Otherwise $f(x) > 0$.*

*b) $f'$ is decreasing, i.e., if $x < y$, then $f'(x) > f'(y)$.*

*Then, given that $f$ satisfies those two conditions, for any pair $x, y \in [0, \infty)$, we have $f(x + y) \leq f(x) + f(y)$.*

*Proof.* we let $g(x) = f(x + y)$ and $p(x) = f(x) + f(y)$, where we regard $y$ as a fixed number. We wish to show that $g(x) \leq p(x)$ or, equivalently, that $0 \leq p(x) - g(x)$. Notice that

$$\frac{d}{dx}(p(x) - g(x)) = p'(x) - g'(x) = f'(x) - f'(x + y) \geq 0 \quad \text{by property b) of } f.$$

Thus, by the first derivative test from the glory days of baby calculus, $p(x) - g(x)$ is increasing for all $x \in [0, \infty)$, attaining its smallest value when $x = 0$. Now,

$$p(0) - g(0) = f(0) + f(y) - f(y) = f(y) - f(y) = 0.$$

Thus, $p(x) - g(x) \geq 0$ for all $x$, and the desired result follows. $\square$

**Theorem 14.** *Let $d \colon M \times M \to [0, \infty)$ be a metric function on $M$ and suppose $f \colon [0, \infty) \to [0, \infty)$ satisfies the two properties discussed in Lemma 1. If $f'(t) > 0$ for all $t \in (0, \infty)$, then $\rho \colon M \times M \to [0, \infty)$ given by $\rho(x, y) = f(d(x, y))$ defines another metric on $M$.*

*Proof.* We have to check that all the properties of metrics are satisfied. Clearly $0 \leq f(d(x, y)) < \infty$ for all $x, y \in M$. Now,

- Suppose $\rho(x, y) = 0$, then $f(d(x, y)) = 0$. By property a) of $f$ (1), this implies that $d(x, y) = 0$, or $x = y$. Obviously $\rho(x, x) = 0$.

- The symmetry $\rho(x, y) = \rho(y, x)$ for all $x, y \in M$ is obvious.

- The triangle inequality is also satisfied:

$$\begin{aligned}
\rho(x, y) &= f(d(x, y)) \\
&\leq f(d(x, z) + d(z, y)) \\
&\leq f(d(x, z)) + f(d(z, y)) \\
&= \rho(x, z) + \rho(z, y) \quad \forall x, y, z \in M.
\end{aligned}$$

Note that the first inequality comes from the assumption that $f$ is increasing and $d(x, y) \leq d(x, z) + d(z, y)$, and the second inequality is a consequence of Lemma 1.

Thus, since $\rho$ satisfies all the required properties, we conclude that it is a metric function, as desired. $\square$

Now, since you're so anxious to work through a few problems on your own, here's something for your entertainment and pleasure:

- Let $\rho, \sigma, \eta \colon M \to \mathbb{R}$. Verify that

$$\rho(x,y) = \sqrt{|x-y|}, \quad \sigma(x,y) = \frac{|x-y|}{1+|x-y|}, \quad \text{and} \quad \eta(x,y) = \log(|x-y|+1),$$

  each define a metric on $M$.

- If $d$ is any metric on $M$ and, as before, $\rho, \sigma, \eta \colon M \to \mathbb{R}$, verify that

$$\rho(x,y) = \sqrt{d(x,y)}, \quad \sigma(x,y) = \frac{d(x,y)}{1+d(x,y)}, \quad \text{and} \quad \eta(x,y) = \log(d(x,y)+1),$$

  are also metrics on $M$.

- Let $\rho, \sigma \colon M \to \mathbb{R}$.

  - Is $\rho(x,y) = \sqrt{\log\left(|x^3 - y^3| + 1\right)}$ a metric function on $M$?
  - How about

$$\sigma(x,y) = \frac{\sqrt{\log\left(|x^3 - y^3| + 1\right)}}{1 + \sqrt{\log\left(|x^3 - y^3| + 1\right)}} \; ?$$

## 1.3.1   Normed Vector Spaces

Some of the most important examples of metric spaces in mathematics are the so-called *normed linear spaces* (also known as *normed vector spaces*), which are algebraic structures that satisfy certain axioms (if you are not familiar with vector spaces, I recommend that you read this subsection only after some perusal of linear algebra). A particular example of great importance in analysis is the function space $C^p[a,b]$, which denotes the set of all continuously differentiable (up to order $p$) functions over the domain $[a,b]$ (there is some subtlety with the notion of differentiability at the endpoints $a$ and $b$, but we will not bother with those issues at this stage). This set is endowed with some additional algebraic structure that makes it a normed vector space. Special cases are the set of all *smooth* (or $C^\infty$) functions

$C^\infty[a,b]$, and also the space of all continuous functions $C[a,b]$ (sometimes also denoted as $C^0[a,b]$.

An easy way to build a metric on a vector space is by way of a length function or norm:

**Definition 10.** *A **norm** on a vector space $V$ is a function $\|\cdot\| : V \to [0,\infty)$ satisfying*

   *i)* $\|x\| = 0 \iff x = 0$.

   *ii)* $\|\alpha x\| = |\alpha|\|x\|$ *for any scalar $\alpha$ and any $x \in V$.*

   *iii)* $\|x + y\| \leq \|x\| + \|y\|$ *for all $x, y \in V$.*

*A function $\|\cdot\| : V \to [0,\infty)$ satisfying all of the above properties except i) is called a **pseudo-norm** (or **semi-norm**); that is, a pseudonorm allows nonzero vectors to have zero length. The pair $(V, \|\cdot\|)$, consisting of a vector space $V$ together with a norm on $V$, is called a **normed vector space**. It is easy to see that any norm induces a metric on $V$ by setting $d(x,y) = \|x - y\|$ (we will refer to this particular metric as the usual metric on $(V, \|\cdot\|)$).*

**Example 6.** *a) The absolute value function $|\cdot|$ clearly defines a norm on $\mathbb{R}$.*

   *b) Each of the following defines a norm on $\mathbb{R}^n$. For $x = (x_1, \ldots, x_n) \in \mathbb{R}^n$,*

- $\|x\|_1 = \sum_{i=1}^n |x_i|$.

- $\|x\|_2 = \left(\sum_{i=1}^n |x_i|^2\right)^{1/2}$.

- *As it happens, continuing the process on the previous two examples for $1 \leq p < \infty$, the expression $\|x\|_p = \left(\sum_{i=1}^n |x_i|^p\right)^{1/p}$ defines a norm on $\mathbb{R}^n$.*

- $\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$.

*The first and last expressions are very easy to check while the second takes a bit more work. The function $\|\cdot\|_2$ is often called the **Euclidean norm** and is generally accepted as the norm of choice in $\mathbb{R}^n$ (no wonder it's the norm you've been seeing since kindergarten math!).*

   *c) Each of the following defines a norm on $C[a,b]$:*

- $\|f\|_1 = \int_a^b |f(t)|\, dt.$

- $\|f\|_2 = \left( \int_a^b |f(t)|^2\, dt \right)^{1/2}.$

- $\|f\|_\infty = \max_{a \leq t \leq b} |f(t)|.$

*Again, the second expression is the hardest to check. The last expression is generally taken as "the" norm on $C[a, b]$.*

   *d) If $(V, \|\cdot\|)$ is a normed vector space, and if $W$ is a linear subspace of $V$, then $W$ is also normed by $\|\cdot\|$. That is, the restriction of $\|\cdot\|$ to $W$ defines a norm on $W$.*

   *e) We now present the sequence space analogues of the "scale" of norms on $\mathbb{R}^n$ given in part b) above. For $1 \leq p < \infty$, we define $\ell_p$ to be the collection of all real sequences $x = \{x_n\}$ for which $\sum_{n=1}^\infty |x_n|^p < \infty$ (by "real sequences" of course we mean $x_i \in \mathbb{R}$ for all i). In the case when $p = \infty$, we define $\ell_\infty$ to be the collection of all bounded real sequences (the definition is cooked up so that we use the $\sup$ norm on such sequences, as you will see now).*

   *It can be shown that each $\ell_p$ is a vector space under coordinatewise addition and scalar multiplication. Moreover, the expression*

$$\|x\|_p = \begin{cases} \left( \sum |x_i|^p \right)^{1/p} & \text{if } 1 \leq p < \infty, \\ \sup_{n \in \mathbb{N}} |x_n| & \text{if } p = \infty. \end{cases}$$

*defines a norm on $\ell_p$. The cases $p = 1$ and $p = \infty$ are not hard to check (do it!). We will verify the other results shortly.* ☘

**Lemma 2 (The Cauchy-Schwarz Inequality).** *For any $x, y \in \ell_2$, we have $\sum_{i=1}^n |x_i y_i| \leq \|x\|_2 \|y\|_2$.*

*Proof.* To simplify notation, we write $\langle x, y \rangle = \sum x_i y_i$. We first consider the case where $x, y \in \mathbb{R}^n$ (that is, $x_i = 0 = y_i$ for all $i > n$). In this case, $\langle x, y \rangle$ is the usual inner product on Euclidean space, i.e., the dot product). Also notice that we may suppose that $x, y \neq 0$ (there is nothing to show if either is 0). Now let $t \in \mathbb{R}$ and consider

$$0 \leq \|x + ty\|_2^2 = \langle x + ty, x + ty \rangle = \|x\|_2^2 + 2t \langle x, y \rangle + t^2 \|y\|_2^2.$$

Since this (nontrivial) quadratic in $t$ is always nonnegative, it must have a nonpositive discriminant. Thus,

$$(2\langle x, y \rangle)^2 - 4\|x\|_2^2 \|y\|_2^2 \le 0 \implies |\langle x, y \rangle| \le \|x\|_2 \|y\|_2.$$

That is,

$$\left| \sum_{i=1}^{n} x_i y_i \right| \le \|x\|_2 \|y\|_2. \tag{1.2}$$

This is not exactly what we were looking for, but it actually implies the stronger inequality in the statement of the lemma. Why? Because the inequality (1.2) must also hold for the vectors $(|x_1|, |x_2|, \ldots, |x_n|)$ and $(|y_1|, |y_2|, \ldots, |y_n|)$. That is,

$$\sum_{i=1}^{n} |x_i y_i| \le \| (|x_1|, \ldots, |x_n|) \| \, \| (|y_1|, \ldots, |y_n|) \| = \|x\|_2 \|y\|_2.$$

Finally, let $x, y \in \ell_2$. Then, for each $n$ we have

$$\sum_{i=1}^{n} |x_i y_i| \le \left( \sum_{i=1}^{n} |x_i|^2 \right)^{1/2} \left( \sum_{i=1}^{n} |y_i|^2 \right)^{1/2} \le \|x\|_2 \|y\|_2.$$

Thus, $\sum_{i=1}^{\infty} x_i y_i$ must be absolutely convergent and satisfy $\sum_{i=1}^{n} |x_i y_i| \le \|x\|_2 \|y\|_2$.    $\square$

Now we are ready to prove the triangle inequality for the $\ell_2$ norm.

**Theorem 15 (Minkowski's Inequality).** *If $x, y \in \ell_2$, then $x + y \in \ell_2$. Moreover, $\|x + y\|_2 \le \|x\|_2 + \|y\|_2$.*

*Proof.* It follows from the Cauchy-Schwarz inequality that, for each $n$, we have

$$
\begin{aligned}
\sum_{i=1}^{n} |x_i + y_i|^2 &= \sum_{i=1}^{n} |x_i|^2 + 2 \sum_{i=1}^{n} x_i y_i + \sum_{i=1}^{n} |y_i|^2 \\
&\le \|x\|_2^2 + 2\|x\|_2 \|y\|_2 + \|y\|_2^2 \\
&= \left( \|x\|_2 + \|y\|_2 \right)^2.
\end{aligned}
$$

Thus, since $n$ is arbitrary, we have $x + y \in \ell_2$ and $\|x + y\|_2 \le \|x\|_2 + \|y\|_2$.    $\square$

We now extend the result of Theorem 15 to general spaces $\ell_p$ for $1 < p < \infty$. Just as in the case of $\ell_2$, several facts are trivial. For example, it is clear that $\|x\|_p = 0$ implies that $x = 0$, and it is easy to see that $\|\alpha x\|_p = |\alpha|\|x\|_p$ for any scalar $\alpha$. Thus we only need to focus our attention on the triangle inequality. We begin with a few classical inequalities that are of interest in their own right. The first shows that $\ell_p$ is at least a vector space.

**Lemma 3.** *Let $1 < p < \infty$ and let $a, b \geq 0$. Then, $(a + b)^p \leq 2^p(a^p + b^p)$. Consequently, $x + y \in \ell_p$ whenever $x, y \in \ell_p$.*

*Proof.* First note that $(a + b)^p \leq (2\max\{a, b\})^p = 2^p \max\{a^p, b^p\} \leq 2^p(a^p + b^p)$. Thus, if $x, y \in \ell_p$, then

$$\sum_{n=1}^{\infty} |x_n + y_n|^p \leq 2^p \sum_{n=1}^{\infty} |x_n|^p + 2^p \sum_{n=1}^{\infty} |y_n|^p < \infty. \qquad \square$$

**Lemma 4 (Young's Inequality).** *Let $1 < p < \infty$ and let $q$ be defined by*

$$\frac{1}{p} + \frac{1}{q} = 1.$$

*Then, for any $a, b \geq 0$, we have*

$$ab \leq \frac{a^p}{p} + \frac{b^q}{q},$$

*with equality occurring if and only if $a^{p-1} = b$.*

*Proof.* Since the inequality trivially holds if either $a$ or $b$ is 0, we may suppose $a, b > 0$. Since $1/p + 1/q = 1$, we see that

$$p\left(\frac{1}{p} + \frac{1}{q}\right) = p \implies 1 + \frac{p}{q} = p.$$

In particular, $p/q = p - 1$. Similarly,

$$q\left(\frac{1}{p} + \frac{1}{q}\right) = q \implies \frac{q}{p} = q - 1.$$

Thus,

$$\frac{1}{p-1} = \frac{1}{p/q} = \frac{q}{p} = q - 1.$$

Also notice that $q = 1/(p-1) + 1$, implying that, just like $p$, $q$ is also in $(1, \infty)$. Thus, the functions $f(x) = x^{p-1}$ and $g(x) = x^{q-1}$ are inverses for $x \geq 0$. The proof of the inequality follows from a comparison of areas:



The area of the rectangle with sides of lengths $a$ and $b$ is at most the sum of the areas under the graphs of the functions $y = x^{p-1}$ for $0 \leq x \leq a$ and $x = y^{q-1}$ for $0 \leq y \leq b$. That is,

$$ab \leq \int_0^a x^{p-1} \mathrm{d}x + \int_0^b y^{q-1} \mathrm{d}y = \frac{a^p}{p} + \frac{b^q}{q}.$$

Clearly, equality can occur only if $a^{p-1} = b$. $\qquad\square$

When $p = q = 2$, Young's inequality reduces to the *arithmetic-geometric mean inequality* (although it is usually stated in the form $\sqrt{ab} \leq (a+b)/2$). Young's inequality will supply the extension of the Cauchy-Schwarz inequality that we need.

**Lemma 5 (Hölder's Inequality).** *Let $1 < p < \infty$ and let $q$ be defined by $1/p + 1/q = 1$. Given $x \in \ell_p$ and $y \in \ell_q$, we have*

$$\sum_{i=1}^{\infty} |x_i y_i| \leq \|x\|_p \|y\|_q.$$

*Proof.* We may suppose that $\|x\|_p > 0$ and $\|y\|_q > 0$ (since otherwise there is nothing to show!). Now, for $n \geq 1$, we use Young's inequality to estimate:

$$\sum_{i=1}^{n} \left| \frac{x_i y_i}{\|x\|_p \|y\|_q} \right| \leq \frac{1}{p} \sum_{i=1}^{n} \left| \frac{x_i}{\|x\|_p} \right|^p + \frac{1}{q} \sum_{i=1}^{n} \left| \frac{y_i}{\|y\|_q} \right|^q \leq \frac{1}{p} + \frac{1}{q} = 1.$$

Thus, we have $\sum_{i=1}^{n} |x_i y_i| \leq \|x\|_p \|y\|_q$ for any $n \geq 1$, and the result follows. $\qquad\square$

**Remark:** Our proof of the triangle inequality will be made easier if we first isolate one of the key calculations. Notice that if $x \in \ell_p$, then the sequence $\{|x_n|^{p-1}\}_{n=1}^{\infty} \in \ell_q$, because $(p-1)q = p$. Moreover, we have

$$\|(|x_n|^{p-1})\|_q = \left( \sum_{i=1}^{\infty} |x_i|^p \right)^{1/q} = \|x\|_p^{p-1}.$$

**Theorem 16 (Minkowski's (General) Inequality).** *Let $1 < p < \infty$. If $x, y \in \ell_p$, then $x + y \in \ell_p$. Moreover, $\|x + y\|_p \leq \|x\|_p + \|y\|_p$.*

*Proof.* We have already shown in Lemma 3 that $x + y \in \ell_p$. To prove the triangle inequality, we once again let $q$ be such that $1/p + 1/q = 1$, and we now use Hölder's inequality to estimate:

$$\sum_{i=1}^{\infty} |x_i + y_i|^p = \sum_{i=1}^{\infty} |x_i + y_i| \cdot |x_i + y_i|^{p-1}$$

$$\leq \sum_{i=1}^{\infty} |x_i| \cdot |x_i + y_i|^{p-1} + \sum_{i=1}^{\infty} |y_i| \cdot |x_i + y_i|^{p-1}$$

$$\leq \|x\|_p \cdot \|(|x_n + y_n|^{p-1})\|_q + \|y\|_p \cdot \|(|x_n + y_n|^{p-1})\|_q$$

$$= \|x + y\|_p^{p-1} \left( \|x\|_p + \|y\|_p \right).$$

That is, $\|x + y\|_p^p \leq \|x + y\|_p^{p-1} \left( \|x\|_p + \|y\|_p \right)$, and the triangle inequality follows. □

### 1.3.2   Limits in Metric Spaces

Now that we have generalized the notion of distance, we are ready to define the notion of limits in abstract metric spaces. Throughout this subsection, unless otherwise specified, we will assume that we are always dealing with a generic metric space $(M, d)$.

**Definition 11.** *Given $x \in M$ and $r > 0$, the set $\mathbb{B}_r(x) = \{y \in M \mid d(x, y) < r\}$ is called the* ***open ball about $x$ of radius $r$***. *(If we want to specify that we are using a given metric $d$, then we write $\mathbb{B}_r^d(x)$.) We may occasionally refer to the set $\overline{\mathbb{B}}_r(x) = \{y \in M \mid d(x, y) \leq r\}$ as the* ***closed ball about $x$ of radius $r$***.

**Example 7.** *a) In $\mathbb{R}$ we have $\mathbb{B}_r(x) = (x - r, x + r)$, the open interval of radius r about x. Similarly, $\overline{\mathbb{B}}_r(x) = [x - r, x + r]$ is the closed interval of radius r about x.*

*b) In $\mathbb{R}^2$ the set $\mathbb{B}_r(x)$ is the open disk of radius r centered at x. The appeareance of $\mathbb{B}_r(x)$ in fact depends on the metric at hand, as we now illustrate.*

- *If d is generated by the norm $\|\cdot\|_1$, then $\mathbb{B}_r(x)$ will look like a square of diameter 2r centered at x (see* Figure 1.2, *where the boundary of $\mathbb{B}_r(x)$ is the "circle" labeled* 1*).*

- *If d is generated by the norm $\|\cdot\|_2$, then $\mathbb{B}_r(x)$ will look like a disk of radius r centered at x (see* Figure 1.2, *where the boundary of $\mathbb{B}_r(x)$ is the (true!) circle labeled* 2*).*

- *If d is generated by the norm $\|\cdot\|_p$, with $1 < p < \infty$, then $\mathbb{B}_r(x)$ will look like a brick with rounded corners. As p gets larger, the brick will assume the appearance of a regular square. That is, if d is generated by the sup norm $\|\cdot\|_\infty$, then $(\mathbb{B}_r(x)$ will look like a square with diameter $2\sqrt{2}r$ centered at x (see* Figure 1.2, *where the boundary of $\mathbb{B}_r(x)$ is the "circle" labeled* $\infty$*).*



Figure 1.2: The shape of the "circle" varies according to the metric used.

*c) In a discrete space $(M, d)$, we have $\mathbb{B}_1^d(x) = \{x\}$ and $\mathbb{B}_2^d(x) = M$.*

*d) In a normed vector space $(V, \|\cdot\|)$, the balls centered at 0 play a special role. In this setting, $\mathbb{B}_r(x) = x + \mathbb{B}_r(0) = \{x + y : \|y\| < r\}$ and $\mathbb{B}_r(0) = r\mathbb{B}_1(0) = \{rx : \|x\| < 1\}$ (for your pure entertainment and pleasure, I leave it to you to check that these statements hold!).*

A subset $A$ of $M$ is said to be ***bounded*** if it is contained in some ball, that is, if $A \subset \mathbb{B}_r(x)$ for some $x \in M$ and some $r > 0$. As it turns out, a subset $A \subseteq M$ is bounded if and only if for any $x \in M$, we have $\sup_{a \in A} d(x, a) < \infty$. Related to this is the ***diameter*** of $A$, which we define as $\operatorname{diam} A = \sup\{d(a, b) \mid a, b \in A\}$.

A ***neighborhood of*** $x$ is any set containing an open ball about $x$ (intuitively, we should think of a neighborhood of $x$ as a "thick" set of points near $x$). We say that a sequence of points $\{x_n\}$ in $M$ ***converges*** to a point $x \in M$ if $d(x_n, x) \to 0$. Now, since this definition is stated in terms of the sequence of real numbers $\{d(x_n, x)\}_{n=1}^{\infty}$, we can easily derive the following equivalent reformulations:

$$\begin{cases} \{x_n\} \text{ converges to } x \text{ if and only if, given any } \varepsilon > 0, \text{ there is} \\ \text{an integer } N \geq 1 \text{ such that } d(x_n, x) < \varepsilon \text{ whenever } n \geq N, \end{cases}$$

or

$$\begin{cases} \{x_n\} \text{ converges to } x \text{ if and only if, given any } \varepsilon > 0, \text{ there is} \\ \text{an integer } N \geq 1 \text{ such that } \{x_n \mid n \geq N \subset \mathbb{B}_\varepsilon(x). \end{cases}$$

If we have that $\{x_n \mid n \geq N\} \subset A$ for some $N$, we say that the sequence $\{x_n\}$ is ***eventually in*** $A$. Thus, our last formulation can be written

$$\begin{cases} \{x_n\} \text{ converges to } x \text{ if and only if, given any } \varepsilon > 0, \\ \text{the sequence } \{x_n\} \text{ is eventually in } \mathbb{B}_\varepsilon(x) \end{cases}$$

or, equally,

$$\begin{cases} \{x_n\} \text{ converges to } x \text{ if and only if it is} \\ \text{eventually in every neighborhood of } x. \end{cases}$$

You can use any of the above expressions to describe the phenomenon of convergence of a series; which formulation you use is a matter of taste, although in some cases one version is better than the others. For instance, note that on the very final version, there is no mention whatsoever of $N$'s and $\varepsilon$'s. Anyhow, just as with real sequences, we usually settle for the shorthand $x_n \to x$ in place of the phrase "$\{x_n\}$ converges to $x$." On occasion we will want to display the set $M$, or the metric $d$ (or both); to this end we may occasionally write $x_n \xrightarrow{d} x$.

**Definition 12.** *A sequence $\{x_n\}$ is said to be a **Cauchy sequence** (also, **clustering sequence**) in $(M, d)$ if, given any $\varepsilon > 0$, there is an integer $N \geq 1$ such that $d(x_m, x_n) < \varepsilon$ whenever $m, n \geq N$. Alternatively, we can say that $\{x_n\}$ is Cauchy if and only if, given $\varepsilon > 0$, there is an integer $N \geq 1$ such that $\operatorname{diam}\{x_n \mid n \geq N\} < \varepsilon$.*

**Definition 13.** *Let $(M, d)$ be a metric space. A sequence $\{x_n\}_{n=1}^{\infty} \subset M$ is said to **converge in** $M$ if there is some $x \in M$ such that, for every $\varepsilon > 0$, there is an integer $N > 0$ that satisfies $d(x_n, x) < \varepsilon$ whenever $n \geq N$.*

We will show shortly that every convergent sequence is Cauchy (cf., Proposition 5). The converse however, is not in general true, i.e., Cauchy sequences need not necessarily converge. For a quick example, consider the sequence $\{1/n\}$ living in the space $M = (0, 1]$ under its usual absolute value metric. Then $\{1/n\}$ is indeed Cauchy, but it does not converge to any point in $M$ (notice that it converges to 0 but $0 \notin M$; see Example 8 below). Notice too that $\{1/n\}$ is a bounded sequence with no convergent subsequence. We will see later on that a Cauchy sequence that has a convergent subsequence does converge.

**Example 8.** *a) Let $M = [0, \infty)$ be endowed with its usual absolute value metric $|\cdot|$, and let $s_n = \{1/n\}_{n=1}^{\infty}$ be a sequence in $M$. Now we deduce the following:*

- *The sequence $s_n$ is convergent. Recall that vaguely familiar expression from baby-calculus: $\lim_{n \to \infty} 1/n = 0$. We show now that $s_n$ does indeed converge to 0. Since $\mathbb{R}$ is an ordered field, we have that*

$$\left| \frac{1}{n} - 0 \right| = \frac{1}{n} < \varepsilon \iff n > \frac{1}{\varepsilon}.$$

  *It follows from the Archimidean property of $\mathbb{R}$ that $n > 1/\varepsilon$ can be achieved for some sufficiently large integer $N$. Thus, if $n \geq N$, it follows that $1/n < \varepsilon$, as desired.*

- *The sequence $s_n$ is Cauchy. For $\varepsilon > 0$, let $N$ be a positive integer such that if $n \geq N$, we have $|1/n - 0| \leq \varepsilon/2$. Then for $m, n \geq N$,*

$$\left| \frac{1}{m} - \frac{1}{n} \right| \leq \left| \frac{1}{n} - 0 \right| + \left| \frac{1}{m} - 0 \right| < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon.$$

*b) Consider the same sequence $s_n$ as in part a), but this time in $M = (0, \infty)$. Then $s_n$ is still Cauchy under the usual absolute value metric inherited by $M$ from $\mathbb{R}$. Notice however, that $\lim_{n \to \infty} 1/n = 0 \notin (0, \infty) = M$. Thus $s_n$ does not converge in $M$.* ✤

As we saw in this example, the matter of convergence depended entirely on the space where we were operating. Now we show that this is not always the case; there are instances where convergence and/or the Cauchy criterion depend on the metric function instead:

**Example 9.** *a) Let $(M,d)$ be the metric space $[0,\infty)$ endowed with he discrete metric*

$$d(x,y) = \begin{cases} 1 & \text{if } x \neq y, \\ 0 & \text{if } x = y. \end{cases}$$

*Then the sequence $s_n = \{1/n\}_{n=1}^{\infty}$ does not converge, since for any $x \in M$, $\mathbb{B}_1^d(x) = \{y \in M \mid d(x,y) < 1\} = \{x\}$. In other words, $s_n$ fails to cluster around $x$.*

    *b) Let $M = (0,\infty)$ with metric $d$ defined by $d(x,y) = |1/x - 1/y|$. Then our good old friend $s_n = \{1/n\}_{n=1}^{\infty}$ is not Cauchy, since $d(1/n, 1/m) = |n - m| \geq 1$ for $m \neq n$.* ✈

Now we present a very important result. The following proposition is reassuring; it tells us that when we go somewhere, we will arrive to one place and one place only. It would have been rather confusing if we had arrived at different places at once, don't you think?!

**Proposition 4.** *Limits are unique. That is, if $x_n \xrightarrow{d} x$ and $x_n \xrightarrow{d} y$, then $x = y$.*

*Proof.* We will show that $d(x,y) = 0$ by proving that the distance between any two points $x$ and $y$ become arbitrarily small; i.e., $d(x,y) < \varepsilon$ for any given $\varepsilon > 0$. Since $x_n \xrightarrow{d} y$, there is some $N > 0$ such that $d(x_n, y) < \varepsilon/2$ whenever $n \geq N$. Similarly, since $x_n \xrightarrow{d} x$, there is some $M > 0$ such that $d(x_n, x) < \varepsilon/2$ whenever $n \geq M$. Letting $k = \max\{M, N\}$, we see that

$$d(x,y) \leq d(x, x_n) + d(x_n, y) < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon \quad \text{whenever } n \geq k. \qquad \square$$

As promised in the comments preceding Example 8, we now show that every convergent sequence is also Cauchy. In addition, we also show that Cauchy sequences are always bounded:

**Proposition 5.** *Consider the sequence $\mathfrak{C}_n = \{x_n \mid n \geq 1\}$.*

a) *If $\mathfrak{C}_n$ is convergent, it is Cauchy.*

b) *If $\mathfrak{C}_n$ is Cauchy, it is bounded.*

*Proof of a).* Suppose $\mathfrak{C}_n$ converges to a limit $x$, i.e., $x_n \overset{d}{\to} x$. Then, for any $\varepsilon > 0$, there is a positive integer $N$ such that $d(x_n, x) < \varepsilon/2$ whenever $n \geq N$. Now,

$$d(x_n, x_m) \leq d(x_n, x) + d(x, x_m) < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon \quad \text{whenever } n, m \geq N.$$

Thus, $\mathfrak{C}_n$ is Cauchy, as desired. $\qquad\square$

*Proof of b).* Now we assume that $\mathfrak{C}_n$ is Cauchy in $(M, d)$ and we want to show boundedness, i.e., we want to find $y \in M$ and $r \in \mathbb{R}$ such that $d(x_n, y) \leq r \ \forall n \geq 1$. Let $\varepsilon > 0$. Then, by the assumption that $\mathfrak{C}_n$ is Cauchy, we have that for some $N > 0$, $d(x_n, x_m) < \varepsilon$ whenever $n, m \geq N$. Now, for indexing purposes, we let $y = x_N$ so that we may set

$$r = \sum_{i=1}^{N-1} d(x_i, x_N) + \varepsilon.$$

Observe that $d(x_n, x_N) < \varepsilon$ whenever $n \geq N$ and $d(x_n, x_N) \leq \sum_{i=1}^{N-1} d(x_i, x_N)$ for $n \leq N - 1$. Thus, $d(x_n, x_N) < r \ \forall n \geq 1$, which means that $\mathfrak{C}_n$ is Cauchy, as desired. $\qquad\square$

**Remark:** Note that although every Cauchy sequence is bounded, not every bounded sequence is, in turn, Cauchy. For an easy example, consider $\{n\}_{n=1}^{\infty} \subset \mathbb{R}$ under the discrete metric. This sequence is definitely bounded, but it is not Cauchy. As another example, notice that the sequence $\{(-1)^n\}_{n=1}^{\infty} \subset \mathbb{R}$ is bounded in any metric, as it has a finite range. However this is not a Cauchy sequence either.

## 1.4  Sequences and Series

You know from the glory days of kindergarten calculus that monotone bounded sequences converge, and that any convergent sequence is necessarily bounded. These two facts together raise the question: Does every bounded sequence converge? "Well, of course not!" –yells the angry mathematician. But just how far from convergent is a typical bounded sequence? To answer this, we want to broaden our definition of limits. Let us start by making the following observation:

Let $\{a_n\}_{n=1}^{\infty}$ be a bounded sequence of real numbers, and consider the sequences:

$$t_n = \inf\{a_n, a_{n+1}, a_{n+2}, ...\} \qquad \text{and} \qquad T_n = \sup\{a_n, a_{n+1}, a_{n+2}, ...\}.$$

Then the sequence $\{t_n\}$ increases, $\{T_n\}$ decreases, and $\sup_{k\in\mathbb{N}} a_k \leq t_n \leq T_n \leq \sup_{k\in\mathbb{N}} a_k$ for all $n$ (make sure that you understand this assessment!). Thus we may speak of $\lim_{n\to\infty} t_n$ as the "lower limit" and $\lim_{n\to\infty} T_n$ as the "upper limit" of our original sequence $\{a_n\}$.

These same considerations are meaningful even if we start with an unbounded sequence $\{a_n\}$, although in that case we will have to allow the values $\pm\infty$ for at least some of the $t_n$'s and/or $T_n$'s. That is, if we allow comparisons to $\pm\infty$, then the $t_n$'s will increase and the $T_n$'s will decrease. Of course we will want to use $\sup_{n\in\mathbb{N}} t_n$ and $\inf_{n\in\mathbb{N}} T_n$ in place of $\lim_{n\to\infty} t_n$ and $\lim_{n\to\infty} T_n$, since "sup" and "inf" have more or less obvious extensions to subsets of the extended real number system $[-\infty, \infty]$, whereas "lim" does not. Even so, we are sure to get caught saying something like "$\{t_n\}$ converges to $+\infty$." But we will pay a stiff penalty for too much rigor here; even a simple fact could have a tedious, unnecessarily long description. Therefore, for the remainder of this section we will interpret words such as "limit" and "converges" in this looser sense.

**Definition 14.** *Given any sequence of real numbers $\{a_n\}$, we define*

$$\liminf_{n\to\infty} a_n = \underline{\lim_{n\to\infty}} a_n = \sup_{n\in\mathbb{N}} \inf_{k\geq n} a_k = \sup_{n\geq 1}\{\inf\{a_n, a_{n+1}, a_{n+2}, \dots\}\}$$

*and*

$$\limsup_{n\to\infty} a_n = \overline{\lim_{n\to\infty}} a_n = \inf_{n\in\mathbb{N}} \sup_{k\geq n} a_k = \inf_{n\geq 1}\{\sup\{a_n, a_{n+1}, a_{n+2}, \dots\}\}.$$

*That is,*

$$\liminf_{n\to\infty} a_n = \sup_{n\in\mathbb{N}} t_n \qquad \left(= \lim_{n\to\infty} t_n \ \text{ if } \{a_n\} \text{ is bounded from below}\right)$$

*and*

$$\limsup_{n\to\infty} a_n = \inf_{n\in\mathbb{N}} T_n \qquad \left( = \lim_{n\to\infty} T_n \ \text{if } \{a_n\} \text{ is bounded from above} \right).$$

*The name "lim inf" is short for **limit inferior** (or **limit infimum**) while "lim sup" is short for **limit superior** or **limit supremum**.*

**Example 10.** *a) Let $a_n = 1/n$, so that $t_n = \inf_{k\geq n} a_k = 0$ and $T_n = \sup_{k\geq n} a_k = 1/n$. Clearly, $\lim_{n\to\infty} t_n = \lim_{n\to\infty}(0) = 0$ and $\lim_{n\to\infty} T_n = 0$; thus $\liminf_{n\to\infty} a_n = \limsup_{n\to\infty} a_n = 0$.*

*b) Let $\{a_n\}_{n=1}^{\infty}$ be the sequence $\{1, -1, 2, -2, 3, -3, \dots\}$. Then $\liminf_{n\to\infty} a_n = -\infty < \infty = \limsup_{n\to\infty} a_n$.*

*c) Let $a_n = \frac{(-1)^n}{1+\frac{1}{n}}$. Then $\liminf_{n\to\infty} a_n = -1$ while $\limsup_{n\to\infty} a_n = 1$.*

*d) Let*

$$a_n = \begin{cases} 1 & \text{if } n \text{ is even,} \\ -1/n & \text{if } n \text{ is odd.} \end{cases}$$

*Then $\limsup_{n\to\infty} a_n = 1$ while $\liminf_{n\to\infty} a_n = 0$.*

Note that a sequence $\{a_n\}$ has a limit if and only if $\limsup_{n\to\infty} = \liminf_{n\to\infty}$ (and both are finite). We use analogous definitions when we take a limit along the real numbers. For example, $\limsup_{y\to x} f(y) = \inf_{\delta>0} \sup_{\|y-x\|<\delta} f(y)$.

We shall now compute the limits of some sequences which occur frequently. We will use the binomial theorem, which was briefly discussed on Section §1.1 in our remarks prior to Theorem 4. We will also need the following fact: If $0 \leq x_n \leq s_n$ for some $n \geq N$ (where $N$ is a fixed number), and if $s_n \to 0$, then it is always true that $x_n \to 0$.

**Theorem 17.** *We have the following results of some special sequences:*

*a) If $p > 0$, then $\lim_{n\to\infty} \frac{1}{n^p} = 0$.*

*b) If $p > 0$, then $\lim_{n\to\infty} \sqrt[n]{p} = 1$.*

*c) $\lim_{n\to\infty} \sqrt[n]{n} = 1$.*

*d)* If $p > 0$ and $\alpha \in \mathbb{R}$, then $\lim_{n \to \infty} \frac{n^{\alpha}}{(1+p)^n} = 0$.

*e)* If $\|x\| < 1$, then $\lim_{n \to \infty} x^n = 0$.

*Proof of a).* Simply consider $n > \left(\frac{1}{\varepsilon}\right)^{1/p}$. (Note that the Archimedian property of $\mathbb{R}$ is used here). $\qquad\square$

*Proof of b).* If $p > 1$, put $x_n = \sqrt[n]{p} - 1$. Then $x_n > 0$ and by the binomial theorem we have

$$1 + n x_n \leq (1 + x_n)^n = p \implies 0 < x_n \leq \frac{p-1}{n},$$

so that $x_n \to 0$. If $p = 1$, b) is trivial, and if $0 < p < 1$, the result is obtained by taking reciprocals. $\qquad\square$

*Proof of c).* Put $x_n = \sqrt[n]{n} - 1$. Then $x_n \geq 0$ and by the binomial theorem we have

$$n = (1 + x_n)^n \geq \frac{n(n-1)}{2} x_n^2 \implies 0 \leq x_n \leq \sqrt{\frac{2}{n-1}} \quad \text{(for } n \geq 2\text{)}. \qquad\square$$

*Proof of d).* Let $k > \alpha$ and $k > 0$. For $n > 2k$, observe that

$$n - k + 1 - \frac{n}{2} = \frac{n}{2} - k + 1 > k - k + 1 = 1 > 0 \qquad \text{so that} \quad n - k + 1 > \frac{n}{2}.$$

Now note that

$$(1 + p)^n > \binom{n}{k} p^k = \frac{n(n-1)\cdots(n-k+1)}{k!} p^k > \frac{n^k}{2^k k!} p^k. \quad \text{(you check that this holds!)}$$

Hence we have

$$0 < \frac{n^{\alpha}}{(1+p)^n} < \frac{2^k k!}{p^k} n^{\alpha - k} \quad \text{(for } n > 2k\text{)}.$$

Since $\alpha - k < 0$, it follows that $n^{\alpha - k} \to 0$ by a), and we are done. $\qquad\square$

*Proof of e).* Simply take $\alpha = 0$ in d), and the desired result follows. $\qquad\square$

For the remainder of this section, all sequences and series under consideration will be complex-valued, unless the contrary is explicitly stated. We now proceed with a very important result. The Cauchy criterion (i.e., that every convergent sequence is Cauchy) can be restated in the following form:

**Theorem 18.** *The sum $\sum a_n$ converges if and only if for every $\varepsilon > 0$, there is an integer $N$ such that $\left\| \sum_{k=n}^{m} a_k \right\| \leq \varepsilon$ if $m \geq n \geq N$. (In particular, by taking $m = n$, the above expression becomes $\|a_n\| \leq \varepsilon$. Also notice that if $n = N$ and $m \to \infty$, the expression becomes $\left\| \sum_{n=N}^{\infty} a_n \right\| \leq \varepsilon$.)*

*Proof.* In $\mathbb{R}$ and in $\mathbb{C}$ every Cauchy sequence converges (this is due to the fact that both $\mathbb{R}$ and $\mathbb{C}$ are examples of what it's known as *complete spaces*. We will discuss this property in detail on Subsection §**??**). Thus the sequence $s_n$ of partial sums is convergent if and only if it is Cauchy (this is *always* the case in complete spaces, although not in general, as we have previously seen when we discussed Cauchy sequences earlier). Now, $s_n$ is Cauchy if for any $\varepsilon > 0$ there is some $N$ such that $n, m \geq N$ implies $\|s_n - s_m\| < \varepsilon$. If $m \geq n$, we have

$$\|s_n - s_m\| = \|s_m - s_n\| = \left\| \sum_{k=1}^{m} a_k - \sum_{k=1}^{n} a_k \right\| = \left\| \sum_{k=n+1}^{m} a_k \right\|. \qquad \square$$

**Corollary 4.** *If $\sum a_n$ converges, then $\lim_{n \to \infty} a_n = 0$.*

The condition $a_n \to 0$ is not, however, sufficient to ensure convergence of $\sum a_n$. For instance the series $\sum_{n=1}^{\infty} 1/n$ diverges(!), as we will demonstrate later.

**Corollary 5.** *A series of nonnegative terms converges if and only if its partial sums form a bounded sequence.*

We now turn to a convergence test of a different nature, the so-called **comparison test** from the glory days of baby-calculus:

**Theorem 19.** *a) If $\|a_n\| \leq c_n$ for $n \geq N_0$ (where $N_0$ is some fixed integer), and if $\sum c_n$ converges, then $\sum a_n$ converges.*
*b) If $a_n \geq d_n \geq 0$ for $n \geq N_0$, and if $\sum d_n$ diverges, then $\sum a_n$ diverges. (Note that this part b) applies only to series of nonnegative terms $a_n$.)*

*Proof.* Given $\varepsilon > 0$, there exists $N \geq N_0$ such that $m \geq n \geq N$ implies $\sum_{k=n}^{m} c_k \leq \varepsilon$ by the Cauchy criterion. Hence

$$\left\| \sum_{k=n}^{m} a_k \right\| \leq \sum_{k=n}^{m} \|a_k\| \leq \sum_{k=n}^{m} c_k \leq \varepsilon,$$

and thus a) follows. Next, b) follows from a), for if $\sum a_n$ converges, so must $\sum d_n$.    □

The comparison test is a very useful one. To use it efficiently however, we should become more familiar with a number of series of nonnegative terms whose convergence or divergence is known. We start this discussion with the the geometric series, which is perhaps the simplest (or at least the one you are more familiar with) of all:

**Theorem 20.** *If $0 \leq x < 1$, then $\sum_{n=0}^{\infty} x^n = \frac{1}{1-x}$. If $x \geq 1$, the series diverges.*

*Proof.* If $x \neq 1$, then we have

$$s_n = \sum_{k=0}^{n} x^k = \frac{1 - x^{n+1}}{1 - x}.$$

The result then follows immediately by letting $n \to \infty$. In the case when $x = 1$ on the hand, we get $1 + 1 + \cdots + 1$, which evidently diverges.    □

In many cases which occur in applications, the terms of the series decrease monotonically. The following theorem of Cauchy is therefore of particular interest. The striking feature of the theorem is that a rather "thin" subsequence of $\{a_n\}$ determines the convergence or divergence of $\sum a_n$.

**Theorem 21.** *Suppose $a_1 \geq a_2 \geq \cdots \geq 0$ (i.e., $\{a_n\}$ is a nonnegative monotonically decreasing sequence). Then the series $\sum_{n=1}^{\infty} a_n$ converges if and only if the series $\sum_{k=0}^{\infty} 2^k a_{2^k}$ converges.*

*Proof.* Since the series under consideration has nonnegative terms, it suffices to consider boundedness of the partial sums. Let $s_n = a_1 + a_2 + \cdots + a_n$ and $t_k = a_1 + 2a_2 + \cdots + 2^k a_{2^k}$. Then, for $n < 2^k$, we have

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

so that

$$s_n \leq t_k. \tag{†}$$

On the other hand, for $n > 2^k$ we have

$$s_n \geq a_1 + a_2 + (a_3 + a_4) + \cdots + (a_{2^{k-1}+1} + \cdots + a_{2^k})$$
$$\geq \frac{1}{2}a_1 + a_2 + 2a_4 + \cdots + 2^{k-1}a_{2^k}$$
$$= \frac{1}{2}t_k,$$

so that

$$2s_n \geq t_k. \tag{††}$$

Then by (†) and (††), the sequences $\{s_n\}$ and $\{t_k\}$ are either both bounded or both un-bounded. $\qquad\square$

**Corollary 6.** *The sum* $\sum \frac{1}{n^p}$ *converges if* $p > 1$ *and diverges if* $p \leq 1$.

*Proof.* Consider the case when $p \leq 0$. In this situation we have $\lim_{n\to\infty} 1/n^p = \lim_{n\to\infty} n^{-p} = \infty$ and the series diverges. Now, when $p > 0$ the sequence $1/n^p$ decreases in which case Theorem 21 applies and we are led to the series

$$\sum_{k=0}^{\infty} 2^k \frac{1}{2^{kp}} = \sum_{k=0}^{\infty} 2^{(1-p)k}.$$

Now, $2^{1-p} < 1$ if and only if $1 - p < 0$, and the result follows by comparison with the geometric series $\sum x^k$, where $x = 2^{1-p}$. $\qquad\square$

**Corollary 7.** *If* $p > 1$, *the sum* $\sum \frac{1}{n(\log n)^p}$ *converges. If* $p \leq 1$, *the series diverges.*

*Proof.* The monotonicity of the logarithmic function implies that the sequence $\{\log n\}$ increases, which in turn implies that the sequence $\{1/(n \log n)\}$ decreases, so that we can apply Theorem 21. This leads us to the series

$$\sum_{k=1}^{\infty} 2^k \frac{1}{2^k (\log 2^k)^p} = \sum_{k=1}^{\infty} \frac{1}{(k \log 2)^p} = \frac{1}{(\log 2)^p} \sum_{k=1}^{\infty} \frac{1}{k^p},$$

and the conclusion follows. $\qquad\square$

Note that the procedure presented above may evidently be continued. For instance,

$$\sum_{n=3}^{\infty} \frac{1}{n \, \log n \cdot \log(\log n)} \quad \text{diverges,}$$

whereas

$$\sum_{n=3}^{\infty} \frac{1}{n \, \log n \cdot (\log(\log n))^2} \quad \text{converges.}$$

Let us now recall the natural base $e$, which is given by the series $e = \sum_{n=0}^{\infty} 1/n!$. Note that since

$$\begin{aligned}
s_n &= 1 + 1 + \frac{1}{1 \cdot 2} + \frac{1}{1 \cdot 2 \cdot 3} + \cdots + \frac{1}{1 \cdot 2 \cdots n} \\
&< 1 + 1 + \frac{1}{2} + \frac{1}{2^2} + \cdots + \frac{1}{2^{n-1}} \\
&< 3,
\end{aligned}$$

the series converges and thus $e$ is well defined. In fact, the series converges very rapidly and allows us to compute $e$ with great accuracy. It is of interest to note that $e$ can also be defined by means of another limit process; the proof provides a good illustration of operations with limits:

**Theorem 22.** *The natural base $e$ is also given by the limit $\lim_{n \to \infty} \left(1 + \frac{1}{n}\right)^n$.*

*Proof.* Let $s_n = \sum_{k=0}^{n} 1/k!$ and $t_n = (1 + 1/n)^n$. Clearly the sequence $s_n$ is monotonically increasing. To see that $t_n$ is also monotonically increasing, observe that by the binomial theorem we have

$$t_n = \left(1 + \frac{1}{n}\right)^n > \left(1 + n\frac{1}{n}\right) = 2.$$

In fact, if we take $\alpha \in \mathbb{R}$ such that $\alpha > -1$ and $\alpha \neq 0$, then $(1+\alpha)^n > 1 + n\alpha$. Now we get

$$\frac{t_{n+1}}{t_n} = \frac{\left(1 + \frac{1}{n+1}\right)^{n+1}}{\left(1 + \frac{1}{n}\right)^n}$$

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

$$= \left(1 + \frac{1}{n}\right)\left(\frac{n^2 + 2n}{(n+1)^2}\right)^{n+1}$$

$$= \left(1 + \frac{1}{n}\right)\left(1 - \frac{1}{(n+1)^2}\right)^{n+1}$$

$$> \left(1 + \frac{1}{n}\right)\left(1 - \frac{1}{n+1}\right) = 1.$$

Thus $t_{n+1} > t_n$ and $t_n$ is increasing, as desired. Now by the binomial theorem,

$$t_n = 1 + 1 + \frac{1}{2!}\left(1 - \frac{1}{n}\right) + \frac{1}{3!}\left(1 - \frac{1}{n}\right)\left(1 - \frac{2}{n}\right) + \cdots$$

$$\cdots + \frac{1}{n!}\left(1 - \frac{1}{n}\right)\left(1 - \frac{2}{n}\right)\cdots\left(1 - \frac{n-1}{n}\right)$$

$$\leq s_n < e.$$

Thus, $\{t_n\}_{n=1}^{\infty}$ is also a bounded sequence, and we have

$$\lim_{n\to\infty} t_n = \alpha \leq e. \tag{$\ddagger$}$$

Next, if $n \geq m$, we get

$$t_n = 1 + 1 + \frac{1}{2!}\left(1 - \frac{1}{n}\right) + \cdots + \frac{1}{m!}\left(1 - \frac{1}{n}\right)\left(1 - \frac{2}{n}\right)\cdots\left(1 - \frac{m-1}{n}\right) + \cdots$$

$$\cdots + \frac{1}{n!}\left(1 - \frac{1}{n}\right)\cdots\left(1 - \frac{n-1}{n}\right)$$

$$\geq 1 + 1 + \frac{1}{2!}\left(1 - \frac{1}{n}\right) + \cdots + \frac{1}{m!}\left(1 - \frac{1}{n}\right)\cdots\left(1 - \frac{m-1}{n}\right).$$

Thus,

$$\alpha > t_n \geq 1 + 1 + \frac{1}{2!}\left(1 - \frac{1}{n}\right) + \cdots + \frac{1}{m!}\left(1 - \frac{1}{n}\right)\cdots\left(1 - \frac{m-1}{n}\right).$$

Now let $n \to \infty$, keeping $m$ fixed, so that

$$\alpha > 1 + 1 + \frac{1}{2!} + \cdots + \frac{1}{m!} = s_m.$$

This means that

$$\alpha \geq \lim_{n \to \infty} s_m = e. \tag{‡‡}$$

Thus, combining (‡) and (‡‡), we have that $\alpha = e$, and we are done. □

In case you are wondering how fast the series $\sum_{n=1}^{\infty} \frac{1}{n!}$ converges, we can estimate that as follows: Let $s_n$ be the same partial sum we defined above (i.e., $s_n = \sum_{k=0}^{n} 1/k!$). Then we have

$$e - s_n = \frac{1}{(n+1)!} + \frac{1}{(n+2)!} + \frac{1}{(n+3)!} + \cdots$$

$$< \frac{1}{(n+1)!} \left( 1 + \frac{1}{(n+1)} + \frac{1}{(n+1)^2} + \cdots \right)$$

$$= \frac{1}{n!\, n},$$

so that

$$0 < e - s_n < \frac{1}{n!n}. \tag{1.3}$$

To illustrate this result, consider $s_{10}$, for instance. Note that $s_{10}$ approximates $e$ with an error less than $10^{-7}$. The inequality (1.3) is of theoretical interest as well, since it enables us to prove the irrationality of $e$, which is presented next.

**Theorem 23.** *The natural base $e$ is an irrational number.*

*Proof.* Suppose, to the contrary, that $e$ is rational. Then $e = p/q$, where $p$ and $q$ are positive integers. By inequality (1.3), we have

$$0 < q! \left( e - s_q \right) < \frac{1}{q}. \tag{1.4}$$

By our assumption, $q!\, e$ is an integer. Moreover, since

$$q!\, s_q = q! \left( 1 + 1 + \frac{1}{2!} + \cdots + \frac{1}{q!} \right)$$

is also an integer, we see that $q!\left(e - s_q\right)$ is an integer as well. But then since $q \geq 1$, inequality (1.4) implies the existence of an integer between 0 and 1, which is obviously not possible. Thus we have reached a contradiction and therefore we conclude that $e$ must be irrational. $\qquad\square$

We now present some important convergence/divergence tests that you have already seen from calculus. This time however, since you are already on your path to becoming a badass mathematician and you are no longer satisfied by hand-waving, we are going to prove these results rigorously.

**Theorem 24 (Root Test).** *Given $\sum a_n$, put $\alpha = \limsup_{n \to \infty} \sqrt[n]{\|a_n\|}$. Then,*

    *a) if $\alpha < 1$, $\sum a_n$ converges.*

    *b) if $\alpha > 1$, $\sum a_n$ diverges.*

    *c) if $\alpha = 1$, the test gives no information.*

*Proof of a).* If $\alpha < 1$, we can choose $\beta$ so that $\alpha < \beta < 1$, and an integer $N$ such that $\sqrt[n]{\|a_n\|} < \beta$ for $n \geq N$. That is, $n \geq N$ implies $\|a_n\| < \beta^n$. Since $0 < \beta < 1$, the sum $\sum \beta^n$ converges. Convergence of $\sum a_n$ follows now from the comparison test. $\qquad\square$

*Proof of b).* If $\alpha > 1$, there is a sequence $\{n_k\}$ such that $\sqrt[n_k]{\|a_{n_k}\|} \to \alpha$. Hence $\|a_n\| > 1$ for infinitely many values of $n$, so that the condition $a_n \to 0$, which is necessary for the convergence of $\sum a_n$, does not hold. $\qquad\square$

*Proof of c).* Consider the series $\sum 1/n$ and $\sum 1/n^2$. We can see that $\sqrt[n]{1/n} \to 1$ and $\sqrt[n]{1/n^2} \to 1$ as well, but the former series diverges while the latter converges. Thus $\alpha = 1$ gives no information whatsoever. $\qquad\square$

**Theorem 25 (Ratio Test).** *The series $\sum a_n$*

    *a) converges if $\limsup_{n \to \infty} \left\| \dfrac{a_{n+1}}{a_n} \right\| < 1$.*

    *b) diverges if $\left\| \dfrac{a_{n+1}}{a_n} \right\| \geq 1$ for all $n \geq n_0$, where $n_0$ is some fixed integer.*

*Proof of a).*  If condition a) holds, we can find $\beta < 1$, and an integer $N$ such that $\|a_{n+1}/a_n\| < \beta$ for $n \geq N$. In particular,

$$\|a_{N+1}\| < \beta \|a_{N+1}\|,$$
$$\|a_{N+2}\| < \beta \|a_{N+1}\| < \beta^2 \|a_N\|,$$
$$\dotfill$$
$$\|a_{N+p}\| < \beta^p \|a_N\|.$$

That is, $\|a_n\| < \|a_N\|\beta^{-N}\beta^n$ for $n \geq N$, and a) follows from the comparison test, since $\sum \beta^n$ converges. $\qquad\square$

*Proof of b).*  If $\|a_{n+1}\| \geq \|a_n\|$ for $n \geq n_0$, it is easy to see that the condition $a_n \to 0$ does not hold, and b) follows. $\qquad\square$

**Remark:** Note that knowing that $\lim_{n\to\infty} a_{n+1}/a_n = 1$ implies nothing about the convergence of $\sum a_n$, as the series $\sum 1/n$ and $\sum 1/n^2$ have demonstrated before.

**Example 11.**  *a) Consider the series*

$$\frac{1}{2} + \frac{1}{3} + \frac{1}{2^2} + \frac{1}{3^2} + \frac{1}{2^3} + \frac{1}{3^3} + \frac{1}{2^4} + \frac{1}{3^4} + \cdots,$$

*for which*

$$\liminf_{n\to\infty} \frac{a_{n+1}}{a_n} = \lim_{n\to\infty} \left(\frac{2}{3}\right)^n = 0$$

$$\liminf_{n\to\infty} \sqrt[n]{a_n} = \lim_{n\to\infty} \sqrt[2n]{\frac{1}{3^n}} = \frac{1}{\sqrt{3}}$$

$$\limsup_{n\to\infty} \sqrt[n]{a_n} = \lim_{n\to\infty} \sqrt[2n]{\frac{1}{2^n}} = \frac{1}{\sqrt{2}}$$

$$\limsup_{n\to\infty} \frac{a_{n+1}}{a_n} = \lim_{n\to\infty} \frac{1}{2} \left(\frac{3}{2}\right)^n = \infty.$$

*The root test indicates convergence whereas the ratio test does not apply.*

*b) The same is true for the series*

$$\frac{1}{2} + 1 + \frac{1}{8} + \frac{1}{4} + \frac{1}{32} + \frac{1}{16} + \frac{1}{128} + \frac{1}{64} + \cdots,$$

*where*

$$\liminf_{n \to \infty} \frac{a_{n+1}}{a_n} = \frac{1}{8}$$

$$\limsup_{n \to \infty} \frac{a_{n+1}}{a_n} = 2,$$

*but*

$$\liminf_{n \to \infty} \sqrt[n]{a_n} = \frac{1}{2}.$$

**Theorem 26.** *For any sequence $\{c_n\}$ of positive numbers, we have*

$$\liminf_{n \to \infty} \frac{c_{n+1}}{c_n} \leq \liminf_{n \to \infty} \sqrt[n]{c_n},$$

$$\limsup_{n \to \infty} \sqrt[n]{c_n} \leq \limsup_{n \to \infty} \frac{c_{n+1}}{c_n}.$$

*Proof.* The proof may be found on [Rudin, 1964, p. 68-69]. □

## 1.5 Open and Closed Sets

Before we start our discussion of open and closed sets, let us probe a bit deeper into our discussion of convergent and Cauchy sequences and answer some of the questions that we left lingering in the previous sections.

**Definition 15.** *Given a sequence $\{x_n\}_{n=1}^{\infty}$, consider a sequence $\{n_k\}_{k=1}^{\infty}$ of positive integers, such that the $n_i$'s are strictly increasing, i.e., $n_1 < n_2 < n_3 < \ldots$. Then the sequence $\{x_{n_k}\}_{k=1}^{\infty}$ is called a **subsequence** of $\{x_n\}_{n=1}^{\infty}$.*

**Notation:** Up until this point we have interchangeably used $\{x_n\}$, $\{x_n\}_{n=1}^{\infty}$, and also $\{x_n \mid n \geq 1\}$. Since we mean the exact same thing by all three, let us from now on simply write $\{x_n\}$ to simplify notation. However, we will make an exception when two indices are in the same expression in order to avoid any confusion. For instance, we write $\{x_{n_k}\}_k$ to mean $\{x_{n_k}\}_{k=1}^{\infty}$, which is telling us that the $n$ is fixed and we are letting $k$ vary from 1 all the way to $\infty$.

**Example 12.** *Let* $\{x_n\} = \{1/n\}$. *We now show some examples of subsequences of* $\{x_n\}$:

a) *Suppose* $\{n_k\}_k = \{2k+1\}$. *Then the subsequence* $\{x_{n_k}\}_k$ *is given by*

$$\{x_{n_k}\}_k = \{x_{2k+1}\} = \left\{\frac{1}{2k+1}\right\} = \left\{\frac{1}{3}, \frac{1}{5}, \frac{1}{7}, \dots\right\}.$$

b) *Suppose* $\{n_k\}_k = \{2^k\}$. *Then the subsequence* $\{x_{n_k}\}_k$ *is given by*

$$\{x_{n_k}\}_k = \{x_{2^k}\} = \left\{\frac{1}{2^k}\right\} = \left\{\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \dots\right\}.$$

c) *Suppose* $\{n_k\}_k = \{2, 1, 3, 4, 5, 6, \dots\}$. *Then* $\{x_{n_k}\}_k$ *is given by*

$$\{x_{n_k}\}_k = \left\{\frac{1}{2}, 1, \frac{1}{3}, \frac{1}{4}, \dots\right\}.$$

*Note that* $\{x_{n_k}\}_k$ ***does not*** *match the order of* $\{x_n\}$ *and therefore fails to be a subsequence of* $\{x_n\}$. *Notice that the* $n_i$'s *in this case are not strictly increasing, as was demanded on our definition of a subsequence. For instance, note that* $n_1 = 2 > 1 = n_2$.

d) *Suppose* $\{n_k\}_k = \{4, 2, 8, 6, 12, 10, \dots\}$. *Then* $\{x_{n_k}\}_k$ *is given by*

$$\{x_{n_k}\}_k = \left\{\frac{1}{4}, \frac{1}{2}, \frac{1}{8}, \frac{1}{6}, \frac{1}{12}, \frac{1}{10}, \dots\right\}.$$

*Note that* $\{x_{n_k}\}_k$ ***does not*** *match the order of* $\{x_n\}$ *either and therefore fails to be a subsequence of* $\{x_n\}$. *Notice that also in this case* $\{n_k\}_k$ *is not strictly increasing.* ✒

Subsequences are useful tools that will later help us to describe key concepts such as *completeness* and *compactness* (cf., Subsection §1.8.3). For now however, let us settle down for some simple analysis of the relationship between subsequences, convergence, and Cauchy sequences.

**Proposition 6.** *If $x_n \xrightarrow{d} x$, then $x_{n_k} \xrightarrow{d} x$ for any subsequence $\{x_{n_k}\}_k$ of $\{x_n\}$.*

*Proof.* We must show that for any $\varepsilon > 0$, there exists $K > 0$ such that $d(x_{n_k}, x) < \varepsilon$ whenever $k \geq K$. Since $x_n \xrightarrow{d} x$, we know that $d(x_n, x) < \varepsilon$ whenever $n \geq N$, where $N$ is some fixed large enough number. Now setting $K = N$, notice that $n_k \geq K$. Thus, when $k \geq N$, we have $d(x_{n_k}, x) < \varepsilon$, as desired. $\qquad\square$

Now we prove an important result, as was promised in Subsection §1.3.2:

**Proposition 7.** *A Cauchy sequence with a convergent subsequence converges.*

*Proof.* We let $\{x_n\}$ be Cauchy and suppose that $\{x_{n_k}\}_k$ is a convergent subsequence with $x_{n_k} \xrightarrow{d} x$. We must show that $x_n \xrightarrow{d} x$. For $\varepsilon > 0$, let $N_1$ be such that $d(x_n, x_m) < \varepsilon/2$ whenever $n, m \geq N_1$. Similarly, let $N_2$ be such that $d(x_{n_k}, x) < \varepsilon/2$ whenever $k \geq N_2$. Now we let $N = \max\{N_1, N_2\}$. If $n, m > N$, then

$$d(x_n, x_{n_m}) \leq d(x_n, x_m) < \frac{\varepsilon}{2} \quad \text{and} \quad d(x_{n_m}, x) < \frac{\varepsilon}{2}.$$

Thus,

$$
\begin{aligned}
d(x_n, x) &\leq d(x_n, x_{n_m}) + d(x_{n_m}, x) \\
&\leq d(x_n, x_m) + d(x_{n_m}, x) \\
&< \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon.
\end{aligned}
$$
$\qquad\square$

**Proposition 8.** *Every subsequence of a Cauchy sequence is itself a Cauchy sequence.*

*Proof.* Let $\{x_n\}$ be Cauchy and suppose that $\{x_{n_k}\}_k$ is any subsequence. For $\varepsilon > 0$, there is some $N > 0$ such that $d(x_n, x_m) < \varepsilon$ whenever $n, m \geq)$. Notice that $n_n \geq n$ and $n_m \geq m$. Hence $d(x_{n_n}, x_{n_m}) < \varepsilon$. $\qquad\square$

**Proposition 9.** *If every subsequence of $\{x_n\}$ has a further subsequence that converges to a certain point $x$, then $\{x_n\}$ converges to $x$ as well.*

*Proof.* Assume, to the contrary, that every subsequence of $\{x_n\}$ has a further subsequence which converges to $x$, but $\{x_n\}$ itself does not converge to $x$. If $x_n \nrightarrow x$ (i.e. if $\{x_n\}$ does not converge to $x$), then infinitely many elements of $\{x_n\}$ are further from $x$ than some $\varepsilon > 0$. That is, the set $A = \{x_n \mid d(x_n, x) \geq \varepsilon\}$ is infinite if $\varepsilon$ is small enough. Notice that the elements of $A$ form a subsequence $\{x_{n_k}\}_k$ of $\{x_n\}$. Our assumption dictates that any sub-subsequence $\{x_{n_{k_t}}\}_t$ must converge to $x$. But all the elements of $\{x_{n_{k_t}}\}_t$ are elements of $A$. In other words, $d(x_{n_{k_t}}, x) \geq \varepsilon$ for all $t \in \mathbb{N}$, implying that $x_{n_{k_t}} \nrightarrow x$, which brings us to a contradiction. $\qquad\square$

### 1.5.1   Equivalent Metrics

We have already seen that the metric at hand often determines which sequences are Cauchy and which sequences converge. Later we will see that the convergent sequences in $(M, d)$ in turn determine the *open* and *closed* sets of $(M, d)$ and therefore the *continuous functions* on $(M, d)$. Given any other metric function $\rho$, we have generally no reason to expect the metric spaces $(M, d)$ and $(M, \rho)$ to have the same convergent sequences. In this subsection we would like to say a few words about metric functions that generate the same convergent sequences.

**Definition 16.** *Two metrics d and $\rho$ on a set M are said to be **equivalent metrics** if they generate the same convergent sequences: that is, $d(x_n, x) \to 0 \iff \rho(x_n, x) \to 0$.*

   It might be comforting to know that most metric functions that we consider on $\mathbb{R}$ (or $[0, \infty)$) are equivalent metrics. The following proposition explains why.

**Proposition 10.** *Let d be a metric on M and suppose that $\rho$ is defined by $\rho(x, y) = f(d(x, y))$, where $f \colon [0, \infty) \to [0, \infty)$ satisfies the following three properties:*

   *i) $f(t) \geq 0$ with equality if and only if $t = 0$.*

   *ii) $f'(t) > 0$ for $t \in (0, \infty)$.*

   *iii) $f''(t) < 0$ for $t \in (0, \infty)$.*

*Then d and $\rho$ are equivalent.*

*Proof.* Notice that $f$ is continuous and invertible (cf., Section §1.7). Since $f'(t) > 0$, we may state that $f^{-1}$ is differentiable and therefore continuous. Let $\mathcal{U}$ be an open interval in $\mathbb{R}$, and then we observe that if $g\colon \mathcal{U} \subset \mathbb{R} \to \mathbb{R}$ is continuous at 0, then for any sequence $\{t_n\} \subset \mathcal{U}$ with $t_n \to 0$, we have $g(t_n) \to g(0)$ (this will be discussed later on when we consider continuous functions and identify their properties). Now suppose that $\{x_n\} \subset M$ with $x_n \xrightarrow{d} x$. Then $d(x_n, x) \to 0$. Since $f$ is continuous, we see that $f(d(x_n, x)) \to f(0) = 0$ by setting $t_n = d(x_n, x)$. Thus,

$$d(x_n, x) \to 0 \implies f(d(x_n, x)) = \rho(x_n, x) \to 0.$$

On the other hand, if $\rho(x_n, x) \to 0$, then $d(x_n, x) = f^{-1}(\rho(x_n, x)) \to 0$ because $f^{-1}$ is also continuous at 0. We have thus established that $d(x_n, x) \to 0$ if and only if $\rho(x_n, x) \to 0$. Hence $d$ and $\rho$ are equivalent, as desired. $\qquad\square$

The following are all equivalent metrics on $\mathbb{R}$:

- $d(x, y) = |x - y|$,

- $\rho(x, y) = \sqrt{|x - y|}$,

- $\eta(x, y) = \frac{|x-y|}{1+|x-y|}$,

- $\zeta(x, y) = \log(|x - y| + 1)$,

- $\xi(x, y) = \frac{\sqrt{\log(|x-y|+1)}}{1+\sqrt{\log(|x-y|+1)}}$.

Note that equivalent metrics preserve convergent sequences. Must they also have the same Cauchy sequences? The answer is *NOoooo!!*, as we can verify in the following example.

**Example 13.** *a)* *Let $M = (0, \infty)$. Then $d(x, y) = |x - y|$ and $\rho(x, y) = |1/x - 1/y|$ are equivalent metrics on M that **do not** generate the same Cauchy sequences:*

- *First, we check that d and $\rho$ are indeed equivalent. Observe that the function $f(t) = 1/t$ is continuous on $(0, \infty)$. Notice also that $f^{-1}(t) = f(t)$; that is, f is its own inverse. Now, if*

*$x_n \xrightarrow{d} x$, then $f(x_n) \to f(x)$; in other words, $|f(x_n) - f(x)| \to 0$. Hence $x_n \xrightarrow{d} x$ implies that $|1/x_n - 1/x| \to 0$, or $x_n \xrightarrow{\rho} x$.*

*On the other hand, $x_n \xrightarrow{\rho} x$ implies that $f^{-1}(x_n) \to f^{-1}(x)$. That is,*

$$x_n \xrightarrow{\rho} x \implies |f^{-1}(x_n) - f^{-1}(x)| = \left| \frac{1}{x_n} - \frac{1}{x} \right| \to 0.$$

*This in turn implies that*

$$\left| f\left( \frac{1}{x_n} \right) - f\left( \frac{1}{x} \right) \right| = |x_n - x| \to 0.$$

*Thus we have shown that $d$ and $\rho$ are equivalent, as desired.*

- *Now to see that $d$ and $\rho$ fail to generate the same Cauchy sequences, notice that $\{1/n\}$ is a Cauchy sequence when it is considered under the metric $d(x, y) = |x - y|$. Under the metric $\rho$ however, $\rho(1/n, 1/m) = |n - m| \geq 1$ if $m \neq n$. Thus $\{1/n\}$ is not Cauchy under $\rho$.*

*b) Let $M = \mathbb{R}^n$. Then*

$$d_1(x, y) = \|x - y\|_1, \quad d_2(x, y) = \|x - y\|_2, \quad \text{and} \quad d_\infty(x, y) = \|x - y\|_\infty$$

*are all equivalent metrics on $\mathbb{R}^n$, because*

$$\|x - y\|_\infty \leq \|x - y\|_2 \leq \|x - y\|_1$$

*while*

$$\|x - y\|_1 \leq n\|x - y\|_\infty \quad \text{and} \quad \|x - y\|_1 \leq \sqrt{n}\, \|x - y\|_2.$$

*You should verify that all three metrics $d_1, d_2, d_\infty$ **do** generate the same Cauchy sequences on $\mathbb{R}^n$. ✿*

Given two metric spaces $(M, d)$ and $(N, \rho)$, we can define a metric on the product $M \times N$ in a variety of ways. Our only requirement is that a sequence of pairs $(a_n, x_n)$ in $M \times N$ should converge precisely when both coordinate sequences $\{a_n\}$ and $\{x_n\}$ converge in $(M, d)$ and $(N, \rho)$, respectively. I will leave it to you as a quick-and-dirty exercise to show that each of the following define metrics on $M \times N$ enjoy this property and that, moreover, they are equivalent:

- $d_1((a,x),(b,y)) = d(a,b) + \rho(x,y)$.

- $d_2((a,x),(b,y)) = (d(a,b)^2 + \rho(x,y)^2)^{1/2}$.

- $d_\infty((a,x),(b,y)) = \max\{d(a,b), \rho(x,y)\}$.

From now on we will refer to any of $d_1, d_2, d_\infty$ as "the" metric on $M \times N$, which we call the *product metric*. The fact that any one of them will serve equally well is because they are equivalent metrics (as I told you to check!).

## 1.5.2   Open Sets

**Definition 17.** *A set $\mathcal{U}$ in a metric space $(M,d)$ is called an **open set** if $\mathcal{U}$ contains a neighborhood of each of its points. In other words, $\mathcal{U}$ is an open set if, given $x \in \mathcal{U}$, there is some $\varepsilon > 0$ such that $\mathbb{B}_\varepsilon(x) \subset \mathcal{U}$.*

    **Remark:** Note that this definition of an open set is suitable only for metric spaces, since otherwise "the ball of radius $\varepsilon$ around $x$ (i.e., $\mathbb{B}_\varepsilon(x)$)" wouldn't mean jack. This is the formulation that we will use at this stage since metric spaces is what we have been studying extensively up to this point. The more general definition of openness will be studied when we get to the subject of *topology* in Chapter 2.

**Example 14.** *a) In any metric space, the whole space $M$ is an open set. The empty set $\varnothing$ is also open (by default).*

    *b) In $\mathbb{R}$, any open interval is an open set. Indeed, given $x \in (a,b)$, let $\varepsilon = \min\{x - a, b - x\}$. Then $\varepsilon > 0$ and $(x - \varepsilon, x + \varepsilon) \subset (a,b)$. The cases $(a, \infty)$ and $(-\infty, b)$ are similar. While we're at it, notice that the interval $[0,1)$, for example, is not open in $\mathbb{R}$ because it does not contain an entire neighborhood of $0$. (**Warning!** When we talk about open sets in a general setting, we need to mention the "topology" of the set in question. For instance, in this case we would say that $\mathbb{R}$ has the "usual topology" so that its open sets are the open intervals. We could have instead chosen $\mathbb{R}$ equipped with the so called "lower limit topology," in which the open intervals are still open sets, but so are half-open intervals such as $[0,1)$. If this is confusing the sh\*t out of you, don't freak out just yet; this is something that will be discussed at length in Chapter 3. For now, I just wanted to bring this up so that you are conscious that we are not dealing with the more general setting yet. I will make no more remarks regarding this matter until we get to cover this material in-depth later on.)*

*c)* In a discrete space, $\mathbb{B}_1(x) = \{x\}$ is an open set for any $x$. It follows that every subset of a discrete space is open.

Following up on the very last example above, we see that in a discrete space an open ball is always open. Let us take this up a notch and show that open balls are open sets on any metric space. Hey, I know what you're thinking..."Bruuuuh... of course an **OPEN** ball <u>has to</u> be open..." Well, not so fast cowboy, we need to prove it!



**Proposition 11.** *For any $x \in M$ and any $\varepsilon > 0$, the open ball $\mathbb{B}_\varepsilon(x)$ is an open set.*

*Proof.* We will use Figure 1.3 to follow the argument of the proof. Let $y \in \mathbb{B}_\varepsilon(x)$. Then



Figure 1.3: Open balls are open.
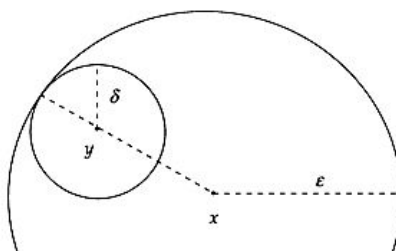
$d(x,y) < \varepsilon$ and thus $\delta = \varepsilon - d(x,y) > 0$. We will show that $\mathbb{B}_\delta(y) \subset \mathbb{B}_\varepsilon(x)$. Indeed, if $d(y,z) < \delta$, then, by the triangle inequality, we have

$$
\begin{aligned}
d(x,z) &\leq d(x,y) + d(y,z) \\
&< d(x,y) + \delta \\
&= d(x,y) + \varepsilon - d(x,y) \\
&= \varepsilon.
\end{aligned}
$$
$\qquad\qquad\square$

Thus we have shown that, as expected, every open ball is open. What's more, if you follow the definition of open sets carefully, you will come to the very wise conclusion that that an open set must actually be a union of open balls: if $\mathcal{U}$ is open, then $\mathcal{U} = \bigcup\{\mathbb{B}_\varepsilon(x) \mid \mathbb{B}_\varepsilon(x) \subset \mathcal{U}\}$. Moreover, any arbitrary union of open balls is again an open set, as the next theorem suggests:

**Theorem 27.** *An arbitrary union of open sets is again open. That is, for any given index A, if $\{\mathcal{U}_\alpha\}_{\alpha \in A}$ is any collection of open sets, then $V = \bigcup_{\alpha \in A} \mathcal{U}_\alpha$ is open.*

*Proof.* If $x \in V$, then $x \in \mathcal{U}_\alpha$ for some $\alpha \in A$. But then, since $\mathcal{U}_\alpha$ is open, it follows that $\mathbb{B}_\varepsilon(x) \subset \mathcal{U}_\alpha \subset V$ for some $\varepsilon > 0$. Since this holds for any given $x \in V$, we are done.     $\square$

So we have shown that open sets behave as nicely as we could possibly want when it comes to unions of such sets. Intersections on the other hand are not so generous; it turns out that we are only allowed to take finitely many of these:

**Theorem 28.** *A finite intersection of open sets is open. That is, if each of $\mathcal{U}_1, \ldots, \mathcal{U}_n$ is open, then so is $W = \mathcal{U}_1 \cap \cdots \cap \mathcal{U}_n$ .*

*Proof.* If $x \in W$, then $x \in \mathcal{U}_i$ for all $i = 1, \ldots, n$. Thus, for each $i$ there is an $\varepsilon_i > 0$ such that $\mathbb{B}_{\varepsilon_i} \subset \mathcal{U}_i$. But then, setting $\varepsilon = \min\{\varepsilon_1, \ldots, \varepsilon_n\} > 0$, we have that $\mathbb{B}_\varepsilon(x) \subset \bigcap_{i=1}^n \mathcal{U}_i = W$. Since this holds for any given $x \in W$, we are done.     $\square$

**Remark:** We have shown in this theorem that we may take any number of intersections (as long as that number is finite) without a problem. If we try to extend this result to infinitely many intersections then, in general, we won't be as lucky. For example, we know that in $\mathbb{R}$, open intervals are open sets; hence $(-1/n, 1/n)$ is open in $\mathbb{R}$ for any $n \in \mathbb{Z}^+$. However $\bigcap_{n=1}^\infty (-1/n, 1/n) = \{0\}$, but $\{0\}$ is not open in $\mathbb{R}$.

The following theorem shows a crucial characteristic of $\mathbb{R}$: every open set in $\mathbb{R}$ can be written as a union of disjoint open intervals. We will see later on that $\mathbb{R}$ does not share this property with its higher-dimensional cousins, i.e., not every open set in $\mathbb{R}^n$, for $n > 1$, can be written as a union of disjoint open balls.

**Theorem 29.** *If $\mathcal{U}$ is an open subset of $\mathbb{R}$, then $\mathcal{U}$ may be written as a countable union of disjoint open intervals. That is, $\mathcal{U} = \bigcup_{n=1}^{\infty} I_n$ , where $I_n = (a_n, b_n)$ (these may be unbounded) and $I_n \cap I_m = \varnothing$ for $n \neq m$.*

*Proof.* We know that $\mathcal{U}$ can be written as a union of open intervals (because each $x \in \mathcal{U}$ is in some open interval $I$ with $I \subset \mathcal{U}$). What we need to show is that $\mathcal{U}$ is a union of disjoint open intervals (I'll leave it to you as a quick exercise to check that such a union must be countable). We first claim that each $x \in \mathcal{U}$ is contained in a maximal open interval $I_x \subset \mathcal{U}$ in the sense that if $x \in I \subset \mathcal{U}$, where $I$ is an open interval, then we must have $I \subset I_x$. Indeed, given $x \in \mathcal{U}$, let

$$a_x = \inf\{a \mid (a, x] \subset \mathcal{U}\} \qquad \text{and} \qquad b_x = \sup\{b \mid [x, b) \subset \mathcal{U}\}.$$

Then, $I_x = (a_x, b_x)$ satisfies $x \in I_x \subset \mathcal{U}$, and $I_x$ is clearly maximal (check this!). Next, notice that for any $x, y \in \mathcal{U}$ we have either $I_x \cap I_y = \varnothing$ or $I_x = I_y$. Why? Because if $I_x \cap I_y \neq \varnothing$, then $I_x \cup I_y$ is an open interval containing both $I_x$ and $I_y$; by maximality we would then have $I_x = I_y$. It thus follows that $\mathcal{U}$ is the union of disjoint (maximal) intervals: $\mathcal{U} = \cup_{x \in \mathcal{U}} I_x$.  □

Now we establish crucial relation between open sets and sequences

**Theorem 30.** *A set $\mathcal{U}$ in $(M, d)$ is open if and only if, whenever a sequence $\{x_n\}_{n=1}^{\infty}$ in M converges to a point $x \in \mathcal{U}$, we have $x_n \in \mathcal{U}$ for all but finitely many n.*

*Proof.* $(\Rightarrow)$ Let $\mathcal{U}$ be an open set containing $x$. We recall that $x_n \to x$ is equivalent to saying that $\{x_n\}$ is eventually in $\mathbb{B}_\varepsilon(x)$ for any $\varepsilon > 0$ (c.f. Section 1.3). But then since $\mathbb{B}_\varepsilon(x) \subseteq \mathcal{U}$, we have that $\{x_n\}$ is eventually in $\mathcal{U}$, which is the same as saying that $x_n \in \mathcal{U}$ for all but finitely many $n$.

$(\Leftarrow)$ We prove this direction by a contrapositive argument, i.e., instead of showing "$\{x_n\}$ is eventually in $\mathcal{U} \implies \mathcal{U}$ is open," we show that "$\mathcal{U}$ not open $\implies \{x_n\}$ is not eventually in $\mathcal{U}$." Thus we start by assuming that $\mathcal{U}$ is not open. Then there is an $x \in \mathcal{U}$ such that $\mathbb{B}_\varepsilon(x) \cap \mathcal{U}^c \neq \varnothing$ for all $\varepsilon > 0$. In particular , for each $n$ there is some $x_n \in \mathbb{B}_{1/n}(x) \cap \mathcal{U}^c$. But then $\{x_n\}_{n=1}^{\infty} \subset \mathcal{U}^c$ and $x_n \to x$, which means that $x_n \notin \mathcal{U}$ for all but finitely many $n$.  □

### 1.5.3  Closed Sets

**Definition 18.** *A set F in a metric space $(M, d)$ is said to be a **closed set** if its complement $F^c = M \smallsetminus F$ is open.*

Well, this definition may not be very enlightening since it depends on the knowledge of open sets, but as we will see below there are intrinsic formulations that do not depend on open sets at all and are actually equivalent to our definition (see Theorem 31 below). In the mean time, let us draw some conclusions about closed sets:

- The entire space $M$ and the empty space $\varnothing$ are always closed. Since, as we saw earlier, $M$ and $\varnothing$ are also both open, this means that it is possible for a set to be both open and closed simultaneously![4] What the frack??!!



Batman's face when he found out that
sets could be both open and closed.

- An arbitrary intersection of closed sets is closed. A finite union of closed sets is closed. (Note that this is the opposite of open sets, which allow arbitrary unions but only finite intersections. This follows by the fact that $(\cap \mathcal{U})^c = \cup \mathcal{U}^c$.)

- Any finite set is closed. Indeed, it is enough to show that a singleton $\{x\}$ is always closed, since any finite set is then given by finite unions of singletons, and finite

---

[4] Such sets that are both open and closed are hilariously called ***clopen sets***.

unions of closed sets are closed. So, let us show that $\{x\}$ is indeed closed: given any $y \in M \smallsetminus \{x\}$, note that $\varepsilon = d(x,y) > 0$, and hence $\mathbb{B}_\varepsilon(y) \subset M \smallsetminus \{x\}$. Thus, since $\{x\}^c = M \smallsetminus \{x\}$ is open, it follows that $\{x\}$ is closed, as desired.

- In $\mathbb{R}$, each of the intervals $[a,b]$, $[a,\infty)$, and $(-\infty,b]$ is closed.

- In a discrete space, every subset is closed. Since, as we previously saw, they are also open, we have that in a discrete space, every subset is clopen!.

- As we have seen above, sets can be open, closed, or both (clopen). We are not done with this lunacy yet: sets can also be neither closed nor open! For instance, $(0,1]$ is neither open nor closed in $\mathbb{R}$!

If your brain is about to explode at this point, don't you worry, you're not alone. As I warned you at the beginning of this chapter, you were in for quite a ride! Just take a coffee break at this point and enjoy watching Hitler struggle with the same stuff that you're struggling with right now: check out *this hilarious video* of Hitler learning basic topology.

Now that you're done laughing at Hitler, let us proceed with our serious business by making an important observation. Let $A$ be a subset of $M$. A point $x \in M$ is called a *limit point* of $A$ if every neighborhood of $x$ contains a point $y \in A$ such that $y \neq x$, that is, if $(\mathbb{B}_\varepsilon(x) \smallsetminus \{x\}) \cap A \neq \varnothing$ for any $\varepsilon > 0$. It is not hard to show that a set $F$ is a closed set if and only if it contains all of its limits points.

Notice that the characterization of closed sets in terms of limit points can readily be translated into a sequential description. To see why, suppose $x$ is a limit point of some set $F$. Then, by definition, $\mathbb{B}_\varepsilon(x) \cap F \neq \varnothing$ for every $\varepsilon > 0$. But this means that for each $1/n$, we can find some $x_n \in \mathbb{B}_{1/n}(x) \cap F$. Thus, the sequence $\{x_n\}_{n=1}^\infty \subset F$ converges to $x$. In other words, any limit point of $F$ is necessarily a point of convergence of some sequence of elements of $F$. This means that $F$ is closed if and only if every sequence $\{x_n\}_{n=1}^\infty$ that consists of elements of $F$ and converges in $M \supset F$, must actually converge to an element of $F$.

We summarize our results in the following theorem:

**Theorem 31.** *Given a set F in $(M,d)$, the following are equivalent:*

   *a)* *F is closed; that is, $F^c = M \smallsetminus F$ is open.*

*b)* If $\mathbb{B}_\varepsilon(x) \cap F \neq \emptyset$ for every $\varepsilon > 0$, then $x \in F$.

*c)* If a sequence $\{x_n\}_{n=1}^\infty \subset F$ converges to some point $x \in M$, then $x \in F$.[5]

*Proof.* (a) $\iff$ b)) This is clear from our observations above and the definition of an open set.

(b) $\implies$ c)) Suppose that $\{x_n\}_{n=1}^\infty \subset F$ and $x_n \to x \in M$. Then $\mathbb{B}_\varepsilon(x)$ contains infinitely many $x_n$ for any $\varepsilon > 0$, and hence $\mathbb{B}_\varepsilon(x) \cap F \neq \emptyset$ for any $\varepsilon > 0$. Thus, $x \in F$ by b).

(c) $\implies$ b)) If $\mathbb{B}_\varepsilon(x) \cap F \neq \emptyset$ for all $\varepsilon > 0$, then for each $n$ there is an $x_n \in \mathbb{B}_{1/n}(x) \cap F$. The sequence $\{x_n\}_{n=1}^\infty$ satisfies $\{x_n\}_{n=1}^\infty \subset F$ and $x_n \to x$. Hence, $x \in F$ by c). $\qquad \square$

As we have seen, some sets are neither open nor closed. However, it is possible to describe the "open part" and the "closure" of a set. We present these concepts briefly here; there will be a more in-depth discussion in Chapter 3.

**Definition 19.** *Given a set $A$ in $(M, d)$, we define the **interior** of $A$, written $\mathrm{Int}\, A$ or $\mathring{A}$, to be the largest open set contained in $A$. That is,*

$$\mathring{A} = \bigcup \{U \mid U \text{ is open and } U \subseteq A\}$$
$$= \bigcup \{\mathbb{B}_\varepsilon(x) \mid \mathbb{B}_\varepsilon(x) \subseteq A \text{ for some } x \in A \text{ and } \varepsilon > 0\}$$
$$= \{x \in A \mid \mathbb{B}_\varepsilon(x) \subseteq A \text{ for some } \varepsilon > 0\}.$$

*On the other hand the **closure** of $A$, usually denoted $\bar{A}$, is defined to be the smallest closed set containing $A$. That is,*

$$\bar{A} = \bigcap \{F \mid F \text{ is closed and } A \subseteq F\}.$$

*Yet other related definitions are the concept of the **exterior** of $A$, denoted by $\mathrm{Ext}\, A$ and given by $\mathrm{Ext}\, A = M \setminus \bar{A}$, and that of the **boundary** of $A$, denoted by $\partial A$ and given by $\partial A = M \setminus (\mathrm{Int}\, A \cup \mathrm{Ext}\, A)$.*

**Proposition 12.** *Given a set $A$, for every $\varepsilon > 0$, we have that $x \in \bar{A}$ if and only if $\mathbb{B}_\varepsilon(x) \cap A \neq \emptyset$.*

---

[5] Note that most authors take this to be the definition of a closed set. That is, a ***closed set*** is usually defined to be a set that contains all of its limit points.

*Proof.* ($\Rightarrow$) Let $x \in \bar{A}$ and assume, to the contrary, that $\mathbb{B}_\varepsilon(x) \cap A = \varnothing$. Then we have that $A$ is a subset of $(\mathbb{B}_\varepsilon(x))^c$, which is a closed set. Thus, $\bar{A} \subset (\mathbb{B}_\varepsilon(x))^c$, since $\bar{A}$ is the smallest closed set containing $A$. But this is a contradiction, because $x \in \bar{A}$ by assumption while $x \notin (\mathbb{B}_\varepsilon(x))^c$.

($\Leftarrow$) If we assume now that $\mathbb{B}_\varepsilon(x) \cap A \neq \varnothing$, then $\mathbb{B}_\varepsilon(x) \cap \bar{A} \neq \varnothing$ (since $A \subset \bar{A}$), and hence $x \in \bar{A}$ (since closed sets contain their limit points). $\qquad \square$

## 1.6 Perfect Sets & Cantor Set

Before we start our discussion of perfect sets, let's present a very important theorem, called *The Nested Interval Theorem*, which is an essential tool that we will be using in the next few topics. Yet before introducing this theorem let's make the following observation:

**Theorem 32.** *A monotone, bounded sequence of real numbers converges in $\mathbb{R}$.*

*Proof.* Let $\{x_n\} \subset \mathbb{R}$ be monotone and bounded. We first suppose that $\{x_n\}$ is increasing (that is, $x_m \leq x_n$ whenever $m < n$). Now, since $\{x_n\}$ is bounded, we may set $x = \sup_{n \in \mathbb{N}}\{x_n\}$ (a real number). We will show that $x = \lim_{n \to \infty} x_n$. Let $\varepsilon > 0$. Since $x - \varepsilon < x = \sup_{n \in \mathbb{N}}\{x_n\}$, we must have $x_N > x - \varepsilon$ for some large enough $N$. But then, for any $n \geq N$, we have $x - \varepsilon < x_N \leq x_n \leq x$. That is, $|x - x_n| < \varepsilon$ for all $n \geq N$. Consequently, $\{x_n\}$ converges and $x = \sup_{n \in \mathbb{N}}\{x_n\} = \lim_{n \to \infty}\{x_n\}$.

Finally, if $\{x_n\}$ is decreasing, consider the increasing sequence $\{-x_n\}$. From the first part of the proof, $\{-x_n\}$ converges to $\sup_{n \in \mathbb{N}}\{-x_n\} = -\inf_{n \in \mathbb{N}}\{x_n\}$. It then follows that $\{x_n\}$ converges to $\inf_{n \in \mathbb{N}}\{x_n\}$. $\qquad \square$

**Theorem 33 (The Nested Interval Theorem).** *If $\{I_n\}$ is a sequence of closed, bounded, nonempty intervals in $\mathbb{R}$ that is decreasing, that is,*

$$I_1 \supseteq I_2 \supseteq I_3 \supseteq \ldots,$$

*Then $\cap_{n=1}^{\infty} I_n \neq \emptyset$. If, in addition, the length of the intervals $I_n$ shrink to 0 (which we denote by* length$(I_n) \to 0$*), then $\cap_{n=1}^{\infty} I_n$ contains precisely one single point.*

**Proof.** Put $I_n = [a_n, b_n]$. Then $I_n \supset I_{n+1}$ means that $a_n \leq a_{n+1} \leq b_{n+1} \leq b_n$ for all $n$. Then we have that

$$a = \lim_{n\to\infty} a_n = \sup_{n\in\mathbb{N}} a_n \qquad \text{and} \qquad b = \lim_{n\to\infty} b_n = \inf_{n\in\mathbb{N}} b_n$$

both exist (as finite real numbers) and satisfy $a \leq b$. Thus, we must have $\cap_{n=1}^{\infty} I_n = [a, b]$. Indeed, if $x \in I_n$ for all $n$, then $a_n \leq x \leq b_n$ for all $n$, and hence $a \leq x \leq b$. Conversely, if $a \leq x \leq b$, then $a_n \leq x \leq b_n$ for all $n$. That is, $x \in I_n$ for all $n$. Finally, if $b_n - a_n =$ length$(I_n) \to 0$, then $a = b$ and so $\cap_{n=1}^{\infty} I_n = \{a\}$. $\qquad\square$

**Example 15.** *a) Note that it is essential that the intervals used in the Nested Interval Theorem be both closed and bounded. Otherwise, for instance, $\cap_{n=1}^{\infty} [n, \infty) = \emptyset$ and $\cap_{n=1}^{\infty} (0, 1, n] = \emptyset$.*

*b) Suppose that $\{I_n\}$ is a sequence of closed intervals with $I_n \supset I_{n+1}$ for all n and with* length$(I_n) \to$ *0 as $n \to \infty$. If $\cap_{n=1}^{\infty} I_n = \{x\}$, then any sequence of points $\{x_n\}$, with $x_n \in I_n$ for all n, must converge to x.*                                                                                                ✺

## 1.6.1   Perfect Sets

**Definition 20.** *A set P is said to be a **perfect set** if it is a closed set and every point of P is a limit point of P. (The empty set $\emptyset$ is declared perfect by default.)*

**Example 16.** *a) The sets $\mathbb{R}$, $(-\infty, a]$, and $[a, \infty)$, as well as any closed and bounded intervals $[a, b]$ (with $a < b$), are all perfect sets.*

*b) The sets $(a, b)$, $[a, b] \cup \{c\}$ (where $b < c$), $\mathbb{Q}$, and $\mathbb{R} \setminus \mathbb{Q}$ are not perfect sets: Note that the sets $(a, b)$, $\mathbb{Q}$, and $\mathbb{R} \setminus \mathbb{Q}$ fail to be closed, even though every point in each of these sets is a limit point of the set. The set $[a, b] \cup \{c\}$ on the other hand, fails to be a perfect set because c is not a limit point of $[a, b] \cup \{c\}$.*

*c) Let $\{x_n\}$ be convergent in $(M, d)$, that is, $x_n \xrightarrow{d} x$ in M. Then the set $\{x_n\} \cup \{x\}$ is not perfect: Although the set is closed, only x is a limit point.*                                                                ✺

Notice that in all of the above listed examples of perfect subsets of $\mathbb{R}$, the perfect sets turned out to be uncountable sets. The next theorem shows that this must always be the case:

**Theorem 34.** *Let $P$ be a perfect subset of $\mathbb{R}$. Then $P$ is uncountable.*

*Proof.* Suppose to the contrary that $P = \{x_1, \ldots, x_n, \ldots\}$ is countable. Let $I_1$ be any closed interval centered at $x_1$ of length $\leq 1$. Then, since $x_1 \in P$ and $P$ is perfect, it follows that $x_1$ is a limit point of $P$. In particular, $(I_1 \setminus \{x_1\}) \cap P \neq \emptyset$.

Now let $n_2$ be the smallest integer for which $x_{n_2} \in (I_1 \setminus \{x_1\}) \cap P$ and let $I_2$ be any closed interval centered at $x_{n_2}$ of length $\leq 1/2$ such that $I_2 \subset I_1$ and $x_1 \notin I_2$. Observe that by the minimality of $n_2$, $x_k \notin I_2$ for any $k < n_2$. Since $x_{n_2} \in P$, it is a limit point of $P$ and therefore $(I_2 \setminus \{x_{n_2}\}) \cap P \neq \emptyset$.

Now let $n_3$ be the smallest integer for which $x_{n_3} \in (I_2 \setminus \{x_{n_2}\}) \cap P$. Set $I_3$ to be any closed interval centered at $x_{n_3}$ of length $\leq 1/3$ such that $I_3 \subset I_2$ and $x_{n_2} \notin I_3$.

Continuing in this fashion, we obtain a nested sequence of closed intervals $I_1 \supset I_2 \supset I_3 \supset \ldots$ such that $\text{length}(I_n) \to 0$ as $n \to \infty$ and $x_k \in I_m$ for all $k < n_m$. By the Nested Interval Theorem, $\cap_{n=1}^{\infty} I_n = \{x\}$ for some $x \in \mathbb{R}$. Notice, however, that $x$ is a limit point of $P$ because it is the limit point of the center points of the intervals $I_n$. Thus, as $P$ is closed, we must have $x \in P$. This is a contradiction: $x$ cannot be any of the $x_m$, since $m \leq n_m$ and $x_m \notin I_{m+1}$. Thus $P$ must be uncountable, as desired. $\qquad\square$

**Example 17.** *a) Although we are still lacking the means to prove it, it can be shown that if $P$ is a nonempty perfect subset of $(M, d)$ in which every Cauchy sequence converges, $P$ must be an uncountable set.*

*b) Let $(M, d)$ be a discrete metric space. Suppose $P \subset M$ is not empty. Then $P$ is not perfect: $\mathbb{B}_1(x) = \{x\}$ for any $x \in P$. Hence $x$ cannot be a limit point of $P$. What went wrong? Every Cauchy sequence in $(M, d)$ must converge in $(M, d)$. Why don't we have perfect sets other than $\emptyset$?*

*c) Let $M = \mathbb{Q}$ under the usual metric of $\mathbb{R}$. Then $P = [0, 1] \cap \mathbb{Q}$ is perfect in $\mathbb{Q}$. Notice, however, that $P \subset \mathbb{Q}$ and must therefore be countable. Does this contradict a)? Absolutely not! Lots of Cauchy sequences of elements of $\mathbb{Q}$ fail to converge in $\mathbb{Q}$.* 🌍

Given these results, it appears that nonempty perfect subsets are rather large. One would expect these sets to occupy a somewhat substantial space on the real number line, for instance. However, reality is stronger than intuition. It seems that perfect subsets of $\mathbb{R}$ can be so constructed as to include almost all of $\mathbb{R}$ and yet be so thin that they fail to contain even a single interval, no matter how small this interval may be. We will see shortly how this materializes when we present the *Cantor set* below, but before we can present our argument, a few definitions will prove useful.

**Definition 21.** *Let $A$ be a subset of a metric space $(M, d)$. If $x \in A$ and $x$ is not a limit point of $A$, then $x$ is called an **isolated point** of $A$. (Note that with this definition, we say that a set is perfect if it is closed and has no isolated points.)*

**Definition 22.** *A set $A$ is said to be **dense** in $(M.d)$ (or, as some authors say, **everywhere dense**) if $\bar{A} = M$. That is, every point of $M$ is either in $A$ or is a limit point of $A$ if $A$ is dense in $M$. Alternatively, we can say that $A \subset M$ is dense in $M$ if $A$ meets every nonempty open subset $W \subset M$, i.e., $A \cap W \neq \varnothing$. On the other hand, a set $A$ is said to be **nowhere dense** in $(M, d)$ if $\text{Int}\,\bar{A} = \varnothing$.*

Note that the intersection of two dense sets need not be dense; it can be empty, as it's the case with $\mathbb{Q}$ and $\mathbb{Q}^c = \mathbb{R} \setminus \mathbb{Q} = \mathbb{I}$. On the other hand if $U, V$ are open dense sets in $M$, then $U \cap V$ is open dense in $M$. For if $W$ is any nonempty open subset of $M$ then $U \cap W$ is a nonempty open subset of $M$ as well, and by denseness of $V$, it is true that $V$ meets $U \cap W$, i.e., $U \cap V \cap W$ is nonempty and $U \cap V$ meets $W$.

**Definition 23.** *A set [space] $X$ is said to be **separable** if it has a countable dense subset [subspace].*

A simple example of a separable space is $\mathbb{R}^n$: Note that $\mathbb{Q}^n$ is a countable dense subset of $\mathbb{R}^n$; thus $\mathbb{R}^n$ is indeed separable.

**Definition 24.** *The countable intersection $G = \cap G_n$ of open dense sets is called a **thick subset** of $M$. Extending our vocabulary in a natural way we say that the complement of a thick set is a **thin subset** (or a **meager subset**).*

**Remark:** Note that a subset $H$ of $M$ is meager if and only if it is a countable union of nowhere dense closed sets, $H = \cup H_n$. Thickness and thinness are "topological" properties.

A thin set is the topological analogue of a zero set (a set whose "outer measure" is zero). These terms in quotation probably don't mean anything to you at this point, but you will see their meaning soon; just hang in there champ!

**Theorem 35 (Baire's Theorem).** *Every thick subset of a complete metric space M is dense in M. A nonempty, complete metric space is not thin: if M is the union of countably many closed sets, then at least one has nonempty interior.*

**Example 18.** *a) The sets $\mathbb{Q}$ and $\mathbb{R} \smallsetminus \mathbb{Q}$ are dense subsets of $\mathbb{R}$.*

*b) Let $A = \{1/n\}$. Then $A$ is nowhere dense in $\mathbb{R}$, because $\bar{A} = A \cup \{0\}$ and $\mathrm{Int}\,\bar{A} = \varnothing$.*

*c) Let $(M,d)$ be discrete. If $A \subseteq M$, then $A$ is both a closed and an open subset of $M$. Thus, $A = \bar{A}$ and $A = \mathrm{Int}\,A$. Thus, if $A$ is nonempty, we have that $A = \bar{A} = \mathrm{Int}\,\bar{A} \neq \varnothing$. Hence, $A$ is not nowhere dense in $M$. Notice, however, that the statement "not nowhere dense" is not equivalent to the phrase "dense" (by now you're probably not even a bit surprised with all this weird stuff, right?). In fact, the only dense subset of $M$ is $M$ itself. Every point of a discrete space is an isolated point.*

*d) Let's elaborate a little bit on the statement mentioned in part c): In general, "not nowhere dense" $\neq$ "dense." For instance, in $\mathbb{R}$, the set $A = (0,1)$ is not dense because $\bar{A} = [0,1] \neq \mathbb{R}$. However, $A$ is not nowhere dense, as $\mathrm{Int}\,\bar{A} = A = (0,1) \neq \varnothing$. The term "nowhere dense" is an unfortunate choice of language that is sadly too common in the literature to avoid. We should think of nowhere dense sets as "thin" sets or sets that are far from having even a single neighborhood.*

*e) While $\mathbb{R}$ is everywhere dense in itself, it is nowhere dense when it is considered as a subset of $\mathbb{R}^2$ (make sure that you understand why this is true).* 🌏

Now we are finally ready to present the argument that we alluded to previous to the above definitions:

**Lemma 6.** *Let $E$ be a closed subset of $(M,d)$ and $E_{iso}$ be the set of all isolated points of $E$. Then $E \smallsetminus E_{iso}$ is a closed subset of $(M,d)$.*

*Proof.* Let $\{x_n\} \subset E \smallsetminus E_{\mathrm{iso}}$ be a sequence that converges to $x \in M$. To show that $E \smallsetminus E_{\mathrm{iso}}$ is closed, we must prove that $x \in E \smallsetminus E_{\mathrm{iso}}$. Notice, however, that since $E$ is closed, $x_n \mapsto x \implies x \in E$. Furthermore, this means that $x$ is a limit point of $E$ so $x$ is not an isolated point, i.e., $x \notin E_{\mathrm{iso}}$. Thus it follows that $x \in E \smallsetminus E_{\mathrm{iso}}$, as desired. $\square$

**Lemma 7.** *Let $E$ be a closed subset of $\mathbb{R}$. Then the set of all isolated points of $E$ (denoted $E_{iso}$ as above) is at most countable.*

*Proof.* Every isolated point of $E$ is contained in an open interval that has no other points of $E$. In other words, if $x \in E_{iso}$, then $x \in I_x$ such that $I_x \cap E = \{x\}$ and $I_x$ is open. Thus $E_{iso} \subset \cup_{x \in E_{iso}} I_x$, where $I_x \cap I_y = \varnothing$ if $x \neq y$. Since each interval contains a rational, this union is countable. This indicates that $E_{iso}$ is also countable as each $I_x$ contains only one point of $E_{iso}$. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\square$

**Theorem 36.** *There exist perfect subsets of $\mathbb{R}$ that contain nearly all of $\mathbb{R}$ and yet fail to have even a single rational number. Such sets must be nowhere dense.*

*Proof.* Let $\varepsilon > 0$. List all the rational numbers in a sequence $\{r_n\}$ and put each $r_n$ at the center of an open interval $I_n$ of length $\varepsilon/2^n$. If $G = \cup_{n=1}^{\infty} I_n$, then $G$ is open and

$$\text{length}(G) \leq \sum_{n=1}^{\infty} \text{length}(I_n) = \sum_{n=1}^{\infty} \frac{\varepsilon}{2^n} = \varepsilon.$$
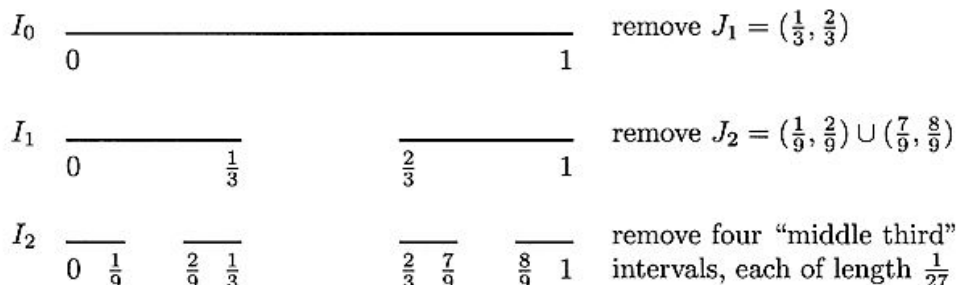
Setting $E = G^c$, we see that $E$ is a closed set whose size must be infinite. That is, $\text{length}(E) = \infty$. Thus, $E$ must be uncountable. Then by Lemma 6, $E \smallsetminus E_{iso}$ is also a closed set and, by Lemma 7, $E \smallsetminus E_{iso}$ must be uncountable (because we delete at most countably many points). Finally, observe that $P = E \smallsetminus E_{iso}$ is a perfect subset of $\mathbb{R}$. $\qquad\qquad\qquad\qquad\qquad\qquad$ $\square$

## 1.6.2 The Cantor Set

Consider the process of successively removing "middle thirds" from the interval $[0,1]$.

We continue this process inductively. At the $n^{th}$ stage we construct $I_n$ from $I_{n-1}$ by removing $2^{n-1}$ disjoint, open, "middle thirds" from $I_{n-1}$, each of length $3^{-n}$; we will call this discarded set $J_n$. Thus, $I_n$ is the union of $2^n$ closed subintervals of $I_{n-1}$, and the complement of $I_n$ in $[0,1]$ is $J_1 \cup \cdots \cup J_n$. The *Cantor set*, which we denote by $\Delta$, is defined as the set of points that still remain at the end of this process; more precisely, $\Delta = \cap_{n=1}^{\infty} I_n$.
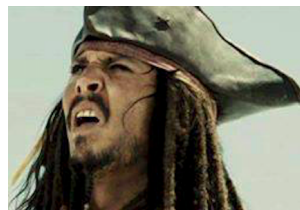
$I_0$  ——————————————————  remove $J_1 = (\frac{1}{3}, \frac{2}{3})$
   0                                                    1

$I_1$  ————    ————  remove $J_2 = (\frac{1}{9}, \frac{2}{9}) \cup (\frac{7}{9}, \frac{8}{9})$
   0          $\frac{1}{3}$        $\frac{2}{3}$        1

$I_2$  ——  ——    ——  ——  remove four "middle third"
   0  $\frac{1}{9}$   $\frac{2}{9}$  $\frac{1}{3}$      $\frac{2}{3}$  $\frac{7}{9}$   $\frac{8}{9}$  1  intervals, each of length $\frac{1}{27}$

It follows from the Nested Interval Theorem that $\Delta \neq \varnothing$, but notice that it is <u>at least</u> countably infinite (in fact we'll show that it's uncountable!). The endpoints of each $I_n$ are in $\Delta$:

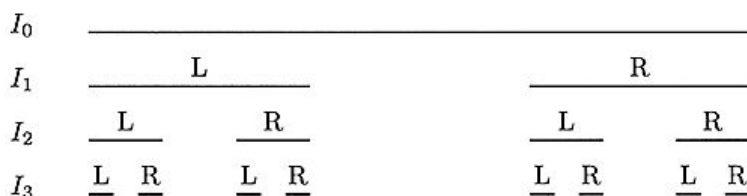$$0, 1, \frac{1}{3}, \frac{2}{3}, \frac{1}{9}, \frac{2}{9}, \cdots \in \Delta.$$

We will refer to these points, which all have the form $a/3^n$ for some integers $a$ and $n$, as the endpoints of $\Delta$.

As we just alluded to, and as we will demonstrate shortly, $\Delta$ turns out to be uncountable! This is another one of those mindblowing results that we have come across thus far. I mean, how on Earth can the Cantor set possibly have the same size as the entire real line??



In what follows we will see two proofs that $\Delta$ is uncountable, the first being somewhat combinatorial. Notice that each subinterval of $I_{n-1}$ results in two subintervals of $I_n$ (after discarding a middle third). We label these two new intervals $L$ and $R$ (for left and right) :

Your face when you were told that the Cantor set is uncountable.

$I_0$  —————————————————————————

$I_1$  ———— L ————        ———— R ————

$I_2$  —— L ——  —— R ——      —— L ——  —— R ——

$I_3$  — L — R —  — L — R —    — L — R —  — L — R —

As we progress down through the levels of the diagram toward the Cantor set (somewhere far below), imagine that we "step down" from one level to the next by repeatedly choosing either a step to the left (landing on an L interval in the next level below) or a step to the right (landing on an R interval). At each stage we are only allowed to step down to a subinterval of the interval we are presently on –jumping across "gaps" is not allowed! Thus, each string of choices, LRLRRLLRLLLR... describes a unique "path" from the top level $I_0$ down to the bottom level $\Delta$.

The Cantor set then is quite literally the "dust at the end of the trail." Said another way, each such LRRLLLR... "path" determines a unique sequence of nested subintervals, one from each level, whose intersection is a single point of $\Delta$. Conversely, each point $x \in \Delta$ lies at the end of exactly one such path, because at any given level there is only one possible subinterval of $I_n$ on our diagram, call it $\widetilde{I}_n$, that contains $x$. The resulting sequence of intervals is clearly nested.

Thus, the Cantor set $\Delta$ is in one-to-one correspondence with the set of all paths, that is, the set of all sequences of L's and R's. Of course, any two choices would have done just as well, so we might also say that $\Delta$ is equivalent to the set of all sequences of 0's and 1's –a set we already know to be uncountable. Here is what this means:

$$\text{card}(\Delta) = 2^{\aleph_0} = \text{card}([0,1]) = \text{card}(\mathbb{R}).$$

Absolutely amazing! The Cantor set is just as "big" as $[0,1]$ (or $\mathbb{R}$) and yet it consists of such a sparse set of points! Unbelievable stuff really...



Before we give our second proof that $\Delta$ is uncountable, let's see why $\Delta$ is also "small" (in a different sense). We will show that $\Delta$ has "measure zero;" that is, the "measure" or

"total length" of all of the intervals in its complement $[0,1] \smallsetminus \Delta$ is 1.[6] In other words, $\Delta$ is not contributing any length whatsoever to the interval $[0,1]$, even though it is composed of uncountably many points!!

Here's the reason for this incredible fact: By induction, the total length of the $2^{n-1}$ disjoint intervals comprising $J_n$ (the set we discard at the $n^{th}$ stage) is $2^{n-1}/3^n$. So the total length of $[0,1] \smallsetminus \Delta$ must be

$$\sum_{n=1}^{\infty} \frac{2^{n-1}}{3^n} = \frac{1}{3} \sum_{n=1}^{\infty} \left(\frac{2}{3}\right)^{n-1} = \frac{1}{3}\frac{1}{1-\frac{2}{3}} = 1.$$

We have discarded everything!? And left uncountably many points behind!? How bizarre! This simultaneous "bigness" and "smallness" is precisely what makes the Cantor set so intriguing.

Our second proof that $\Delta$ is uncountable is based on an equivalent characterization of $\Delta$ in terms of ternary (base 3) decimals. Recall that each $x$ in $[0,1]$ can be written, in possibly more than one way, as $x = 0.a_1 a_2 a_3 \cdots$ (base 3), where each $a_n = 0, 1,$ or 2. This three-way choice for decimal digits (base 3) corresponds to the three-way splitting of intervals that we saw earlier. To see this, let us consider a few specific examples. For instance, the three cases $a_1 = 0, 1,$ or 2 correspond to the three intervals $[0, 1/3], (1/3, 2/3),$ and $[2/3, 1]$:

$$
\begin{array}{c}
I_1 \quad \overset{\displaystyle a_1 = 0}{\rule{3cm}{0.4pt}} \quad \overset{\displaystyle a_1 = 1}{\phantom{xxx}} \quad \overset{\displaystyle a_1 = 2}{\rule{3cm}{0.4pt}} \\
0 \qquad\qquad \tfrac{1}{3} \qquad\qquad \tfrac{2}{3} \qquad\qquad 1
\end{array}
$$

Notice that there is some ambiguity at these endpoints:

$$\frac{1}{3} = 0.1 \text{ (base 3)} = 0.0222\ldots \text{ (base 3)},$$

$$\frac{2}{3} = 0.2 \text{ (base 3)} = 0.1222\ldots \text{ (base 3)},$$

$$1 = 1.0 \text{ (base 3)} = 0.2222\ldots \text{ (base 3)},$$

but each of these ambiguous cases has at least one representation with $a_1$ in the proper range. Now we show the situation for $I_2$ (but this time ignoring the discarded intervals):

---

[6] This concept of "measure" is the main topic of Section 1.9. This is merely a glimpse of what's coming ☺.

| $I_2$ | $a_1 = 0$ and | | | | $a_1 = 2$ and | | | |
|---|---|---|---|---|---|---|---|---|
| | $a_2 = 0$ | | $a_2 = 2$ | | $a_2 = 0$ | | $a_2 = 2$ | |
| | $0$ | $\frac{1}{9}$ | $\frac{2}{9}$ | $\frac{1}{3}$ | $\frac{2}{3}$ | $\frac{7}{9}$ | $\frac{8}{9}$ | $1$ |

Again, some confusion is possible at the endpoints:

$$\frac{1}{9} = 0.01 \text{ (base 3)} = 0.00222\ldots \text{ (base 3)},$$

$$\frac{8}{9} = 0.22 \text{ (base 3)} = 0.21222\ldots \text{ (base 3)}.$$

These examples demonstrate the validity of the following theorem:

**Theorem 37.** *An element $x$ is in the Cantor set if and only if $x$ can be written as $\sum_{n=1}^{\infty} a_n/3^n$, where each $a_n$ is either $0$ or $2$.*

Thus the Cantor set $\Delta$ consists of those points in $[0, 1]$ having some base 3 decimal representation that excludes the digit 1. Knowing this, we can list all sorts of elements of $\Delta$. For example, $1/4 \in \Delta$, because $1/4 = 0.020202\ldots$ (base 3).

## 1.7 Continuity

We start off this section by discussing limits and continuity for real-valued functions. We will then gradually generalize this crucial concept, extending it to more abstract metric spaces. Finally in Chapter 3, we will see continuity in all its glory, when discussed in general (not necessarily metric) topological spaces.

To begin, we let $f$ be a real-valued function defined (at least) for all points in some open interval containing the point $a \in \mathbb{R}$, except possibly at $a$ itself. We will normally refer to such

a set as a ***punctured neighborhood*** of $a$. Given a number $L \in \mathbb{R}$, we write $\lim_{x \to a} f(x) = L$ to mean:

$$\begin{cases} \text{for every } \varepsilon > 0, \text{ there exists some } \delta > 0 \text{ such that } |f(x) - L| < \varepsilon \\ \text{whenever } x \text{ satisfies } 0 < |x - a| < \delta. \end{cases}$$

We then say that $\lim_{x \to a} f(x)$ exists if there is some number $L \in \mathbb{R}$ that satisfies these requirements.

**Theorem 38.** *Let $f$ be a real-valued function defined in some punctured neighborhood of $a \in \mathbb{R}$. Then, the following are equivalent:*

   *i)* *There exists a number L such that $\lim_{x \to a} f(x) = L$ (by the $\varepsilon - \delta$ definition described above).*

   *ii)* *There exists a number L such that $f(x_n) \to L$ whenever $x_n \to a$, where $x_n \neq a$ for all n.*

   *iii)* *$\{f(x_n)\}$ converges (to something) whenever $x_n \to a$, where $x_n \neq a$ for all n.*

   **Remark:** The point to item iii) is that if $\lim_{n \to \infty} f(x_n)$ always exists, then it must actually be independent of the choice of $\{x_n\}$. Indeed, if $x_n \to a$ and $y_n \to a$ as well, then the sequence $x_1, y_1, x_2, y_2, \ldots$ also converges to $a$. This particular phrasing is interesting because it does not refer to $L$. That is, we can test for the existence of a limit without knowing its actual value.

   Now suppose that $f$ is defined in a neighborhood of $a$, this time including the point $a$ itself. We say that $f$ is ***continuous at*** $a$ if $\lim_{x \to a} f(x) = f(a)$. That is, if:

$$\begin{cases} \text{for every } \varepsilon > 0, \text{ there exists some } \delta > 0 \text{ (that depends on } f, a, \text{ and } \varepsilon) \\ \text{such that } |f(x) - f(a)| < \varepsilon \text{ whenever } x \text{ satisfies } |x - a| < \delta. \end{cases}$$

Notice that we replaced $L$ by $f(a)$ and we dropped the requirement that $x \neq a$. You should be already somewhat familiar with this $\varepsilon - \delta$ definition of continuity from your days of kindergarten-calculus, although in those dreadful days all you probably cared about was doing the computations and all the beautiful underlying theory didn't matter much to you. Thankfully you're a grown up person by now and you are dying to learn everything there is to know about continuity and all its ramifications!

   Anyhow, now we can naturally extend Theorem 38 in the obvious way:

**Theorem 39.** *Let $f$ be a real-valued function defined in some neighborhood of $a \in \mathbb{R}$. Then, the following are equivalent:*

   *i)*  *$f$ is continuous at $a$ (by the $\varepsilon - \delta$ definition).*

   *ii)*  *$f(x_n) \to f(a)$ whenever $x_n \to a$.*

   *iii)*  *$\{f(x_n)\}$ converges (to something) whenever $x_n \to a$.*

Notice that we dropped the requirement that $x_n \neq a$; thus if $\lim_{n \to \infty} f(x_n)$ always exists, then it must equal $f(a)$. You might also recall from baby-calculus that there is a standard notation for left and right-hand limits and left and right continuity. So, if we define

$$f(a-) = \lim_{x \to a^-} f(x) \qquad \text{and} \qquad f(a+) = \lim_{x \to a^+} f(x)$$

(provided that these limits exist, of course), then a requirement of continuity of $f$ is that both $f(a-)$ and $f(a+)$ exist, and are both equal to $f(a)$.

This concept of one-sided limits works like a charm on functions defined on $\mathbb{R}$, thanks to the well defined order of $\mathbb{R}$.

However, these one-sided limit do not generalize very well to other spaces, although they are instrumental in allowing us to catalogue different types of discontinuities. For instance, we say that $f$ is **right-continuous at** $a$ if $f(a+)$ exists and equals $f(a)$, and similarly we say that $f$ is **left-continuous at** $a$ if $f(a-)$ exists and equals $f(a)$. We also say that $f$ has a **jump discontinuity at** $a$ if $f(a-)$ and $f(a+)$ both exist but at least one is different from $f(a)$. A function having only jump discontinuities is not that terrible. In particular, monotone functions are rather well behaved:

**Proposition 13.** *Let $f \colon (a, b) \to \mathbb{R}$ be monotone and let $a < c < b$. Then, $f(c-)$ and $f(c+)$ both exist, and so $f$ can have only jump discontinuities.*

*Proof.* We might as well suppose that $f$ is increasing (otherwise, we consider $-f$). In that case, $f(c)$ is an upper bound for $\{f(t) \mid a < t < c\}$ and a lower bound for $\{f(t) \mid c < t < b\}$. All that remains is to check that

$$\sup\{f(t) \mid a < t < c\} = \lim_{x \to c^-} f(x) \qquad \text{and} \qquad \inf\{f(t) \mid c < t < b\} = \lim_{x \to c^+} f(x).$$

We will now sketch the proof of the first of these; I leave it to you as an exercise to check the inf statement:

Given $\varepsilon > 0$, there is some $x_0$ with $a < x_0 < c$ such that $\sup_{t<c} f(t) - \varepsilon < f(x_0) \leq \sup_{t<c} f(t)$. Now let $\delta = c - x_0 > 0$. Then, if $c - \delta < x < c$, we get $x_0 < x < c$, and so $f(x_0) \leq f(x) \leq \sup_{t<c} f(t)$. Thus, $|f(x) - \sup_{t<c} f(t)| < \varepsilon$. $\qquad\square$

**Theorem 40.** *If $f\colon (a,b) \to \mathbb{R}$ is monotone, then $f$ has at most countably many points of discontinuity in $(a,b)$, all of which are jump discontinuities.*

*Proof.* That $f$ has only jump discontinuities follows from Proposition 13. Now we just need to count the points of discontinuity. Let's reflect on the situation. If $f\colon (a,b) \to \mathbb{R}$ is, say, increasing, and if $c \in (a,b)$, then the left and right-hand limits of $f$ at $c$ satisfy $f(c-) \leq f(c) \leq f(c+)$. In particular, $f$ is discontinuous at $c$ if and only if $f(c-) < f(c+)$. Consequently, if $c$ and $d$ are two different points of discontinuity for $f$, then the intervals $(f(c-), f(c+))$ and $(f(d-), f(d+))$ are nonempty and disjoint. Thus,

$$\{(f(c-), f(c+)) \mid c \text{ is a point of discontinuity for } f\}$$

is a collection of nonempty, disjoint open intervals in $\mathbb{R}$, and as we know any such collection must be countable. $\qquad\square$

**Corollary 8.** *If $f\colon [a,b] \to [c,d]$ is both monotone and onto, then $f$ is continuous.*

We can put this corollary to good use. Recall that the Cantor function $f\colon \Delta \to [0,1]$ is monotone and onto. Indeed, if $x \in \Delta$, then $x = 0.2a_1 2a_2 \ldots$ (base 3), where each $a_i = 0$ or 1 and $f(x) = \sum_{i=1}^{\infty} a_i/2^i$. Since $\{a_n\}$ can be any sequence of 0's and 1's, $f$ is clearly onto.

Now we can extend this definition of the Cantor function $f$ to all of $[0,1]$ in an obvious way: We take $f$ to be an appropriate constant on each of the open intervals that make up $[0,1] \setminus \Delta$. For example, we would set

$$f(x) = \begin{cases} f(1/3) = 1/2 & \text{for each } x \in (1/3, 2/3), \\ f(1/9) = 1/4 & \text{for each } x \in (1/9, 2/9). \end{cases}$$

Formally, we define

$$f(x) = \sup\{f(y) \mid y \in \Delta, y \leq x\} \qquad \text{for each } x \in [0,1] \setminus \Delta.$$

The new function $f\colon [0,1] \to [0,1]$ is still increasing and is actually continuous (because it is onto)!! It is called a ***singular function*** because $f' = 0$ at almost every point in $[0,1]$. That is $f' = 0$ on $[0,1] \smallsetminus \Delta$, which is a set of measure 1. We will see more on this so called *Cantor-Lebesgue* function on Subsection §1.11.4.

Theorem 40 has a converse. Given any countable set $\mathcal{D}$ in $\mathbb{R}$, we can construct an increasing function $f\colon \mathbb{R} \to \mathbb{R}$ that is discontinuous precisely at the points of $\mathcal{D}$. Here is a brief sketch:

Let $\mathcal{D} = \{x_1, x_2, \ldots\}$ and let $\{\varepsilon_n\}$ be a sequence of positive numbers with $\sum_{n=1}^{\infty} \varepsilon_n < \infty$. We define $f(x) = \sum_{x_n \leq x} \varepsilon_n$, where the sum is over the set $\{n \mid x_n \leq x\}$ and where $f(x) = 0$ if the set is empty (notice that $0 \leq f(x) \leq \sum_{n=1}^{\infty} \varepsilon_n < \infty$ in any case).

Now, if $x < y$, we have

$$
\begin{aligned}
f(y) = \sum_{x_n \leq y} \varepsilon_n &= \sum_{x_n \leq x} \varepsilon_n + \sum_{x < x_n \leq y} \varepsilon_n \\
&= f(x) + \sum_{x < x_n \leq y} \varepsilon_n \\
&\geq f(x).
\end{aligned}
$$

Thus, $f$ is increasing.

Next we consider this formula in each of the cases $x = x_k$ and $y = x_k$.

Case 1   We have

$$
x = x_k < y \implies f(y) = f(x_k) + \sum_{x_k < x_n \leq y} \varepsilon_n.
$$

We claim that $f(x_k+) = f(x_k)$:

$$
\lim_{y \to x_k^+} \sum_{x_k < x_n \leq y} \varepsilon_n = 0 \qquad \text{because} \qquad \sum_{n=N}^{\infty} \varepsilon_n \to 0 \text{ as } N \to \infty.
$$

Case 2   We have

$$
x < x_k = y \implies f(x_k) = f(x) + \sum_{x < x_n \leq x_k} \varepsilon_n \geq f(x) + \varepsilon_k.
$$

We claim that $f(x_k-) = f(x_k) - \varepsilon_k$; that is,

$$
\lim_{y \to x_k^-} \sum_{x < x_n \leq x_k} \varepsilon_n = \varepsilon_k.
$$

Putting this all together, we have

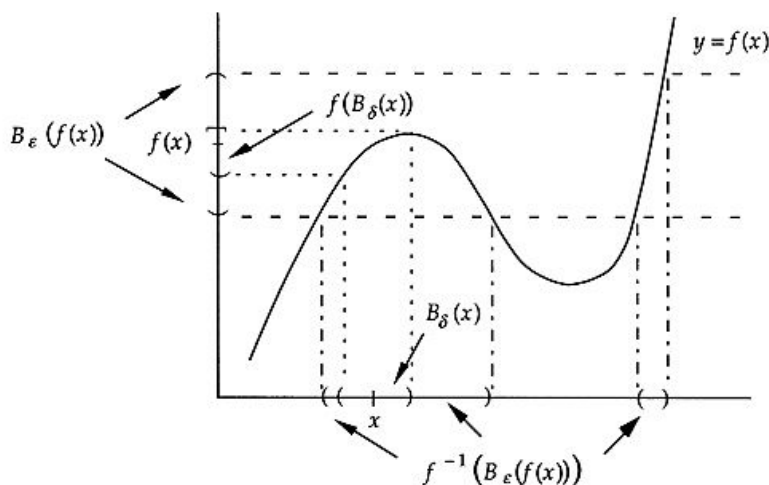$$f(x_k-) + \varepsilon_k = f(x_k) = f(x_k+) \quad \text{and} \quad f(x_k+) - f(x_k-) = \varepsilon_k.$$

The proof that $f$ is continuous at each $x \in \mathbb{R} \setminus \mathcal{D}$ is similar to this procedure.

### 1.7.1 Continuity on Abstract Metric Spaces

Given a mapping $f \colon (M,d) \to (N,\rho)$ (where $M, N$ are arbitrary metric spaces), and given a point $x \in M$, we have at least two plausible definitions for the continuity of $f$ at $x$. Each definition is derived from its obvious counterpart for real-valued functions by replacing absolute values with an appropriate metric. For example, we might say that $f$ is **continuous at** $x$ if $\rho(f(x_n), f(x)) \to 0$ whenever $d(x_n, x) \to 0$. That is, $f$ should send sequences converging to $x$ into sequences converging to $f(x)$. This says that $f$ "commutes" with limits: $f(\lim_{n \to \infty}(x_n)) = \lim_{n \to \infty} f(x_n)$.

Another alternative is to use the familiar $\varepsilon - \delta$ definition from kindergarten-calculus. In this case we would say that $f$ is **continuous at** $x$ if, given any $\varepsilon > 0$, there always exists a $\delta > 0$ such that $\rho(f(x), f(y)) < \varepsilon$ whenever $d(x, y) < \delta$. Written in slightly different terms, this definition requires that $f(\mathbb{B}_\delta^d(x)) \subset \mathbb{B}_\varepsilon^\rho(f(x))$. That is, $f$ maps a sufficiently small neighborhood of $x$ into a given neighborhood of $f(x)$.

We will rewrite this last definition once more, but this time we will use an inverse image. (Use the figure on the right for guidance.) Recall that the *inverse image* of a set $A \subseteq Y$, under a function $f \colon X \to Y$, is defined to be the set $\{x \in X \mid f(x) \in A\}$ and it is usually written $f^{-1}(A)$ (note that while inverse functions are not always defined, the inverse image of any set under any function al-ways makes sense). Stated in terms of an inverse image, our condition reads: $\mathbb{B}_\delta^d(x) \subset$

$f^{-1}(\mathbb{B}_\varepsilon^\rho(f(x)))$. This condition tells us that the inverse image of an open set containing $f(x)$ must still be open near $x$.

Lastly, if $f$ is continuous at every point of $M$, we simply say that $f$ is **continuous on $M$**, or simply that $f$ is continuous. If $f$ has an inverse $f^{-1}$ that is also continuous, then we say that $f$ is a **homeomorphism**. We will postpone any further discussion on this concept until we get to Chapter 3. Homeomorphisms are the key ingredient in topology, as those are the maps that decide which spaces are "equivalent" (topologically speaking). That is, if you can define a homeomorphism between two "topological spaces" $M$ and $N$, then $M$ and $N$ are said to be **homeomorphic** or **topologically equivalent**. That's all you get for now, but there's plenty of entertainment awaiting you in topology so don't fret!

It is crystal clear that all these statements concerning arbitrary open balls can be easily translated into statements concerning arbitrary open sets. Hence, there is a characterization of continuity available that may be stated exclusively in terms of open sets. Of course, any statement concerning open sets probably has a counterpart using closed sets. Moreover, recall that we also characterized open sets and closed set in terms of convergent sequences, and as such we would expect to find a characterization of continuity in terms of convergent sequences as well. Anyhow, it is time to present an extension of Theorem 39 to abstract metric spaces:

**Theorem 41.** *Given $f \colon (M,d) \to (N,\rho)$, the following are equivalent:*

   *i)* *$f$ is continuous on M (by the $\varepsilon - \delta$ definition).*

   *ii)* *For any $x \in M$, if $x_n \to x$ in M, then $f(x_n) \to f(x)$ in N.*

   *iii)* *If E is closed in N, then $f^{-1}(E)$ is closed in M.*

   *iv)* *If V is open in N, then $f^{-1}(V)$ is open in M.*

*Proof.* $(i) \iff ii))$ Suppose that $x_n \xrightarrow{d} x$ . Given $\varepsilon > 0$, let $\delta > 0$ be such that $f\left(B_\delta^d(x)\right) \subset B_\varepsilon^\rho(f(x))$ . Then, since $x_n \xrightarrow{d} x$, we have that $\{x_n\}$ is eventually in $B_\delta^d(x)$. But this implies that $\{f(x_n)\}$ is eventually in $B_\varepsilon^\rho(f(x))$ . Since $\varepsilon$ is arbitrary, this means that $f(x_n) \xrightarrow{\rho} f(x)$. (Compare the result discussed here with the case $f \colon \mathbb{R} \to \mathbb{R}$ that was treated on part i) of Theorem 39.)

(*ii*) $\iff$ *iii*)) Let $E$ be closed in $(N, \rho)$. Given $\{x_n\} \subset f^{-1}(E)$ such that $x_n \xrightarrow{d} x \in M$, we need to show that $x \in f^{-1}(E)$. But $\{x_n\} \subset f^{-1}(E)$ implies that $\{f(x_n)\} \subset (E)$, while $x_n \xrightarrow{d} x \in M$ tells us that $f(x_n) \xrightarrow{\rho} f(x)$ from ii). Thus, since $E$ is closed, we have that $f(x) \in E$ or $x \in f^{-1}(E)$.

(*iii*) $\iff$ *iv*)) This part is rather obvious, since $f^{-1}(A^c) = \left(f^{-1}(A)\right)^c$.

(*iv*) $\iff$ *i*)) Given $x \in M$ and $\varepsilon > 0$, the set $B_\varepsilon^\rho(f(x))$ is open in $(N, \rho)$ and so, by iv), the set $f^{-1}\left(B_\varepsilon^\rho(f(x))\right)$ is open in $(M, d)$. But then $B_\delta^d(x) \subset f^{-1}\left(B_\varepsilon^\rho(f(x))\right)$, for some $\delta > 0$, because $x \in f^{-1}\left(B_\varepsilon^\rho(f(x))\right)$. $\qquad\square$

**Example 19.** *a) Consider the **characteristic function** of* $\mathbb{Q}$, $\chi_{\mathbb{Q}} \colon \mathbb{R} \to \mathbb{R}$, *which is given by*

$$\chi_{\mathbb{Q}}(x) = \begin{cases} 1 & \text{if } x \in \mathbb{Q}, \\ 0 & \text{if } x \notin \mathbb{Q}. \end{cases}$$

*Then* $\chi_{\mathbb{Q}}^{-1}(\mathbb{B}_{1/3}(1)) = \mathbb{Q}$ *and* $\chi_{\mathbb{Q}}^{-1}(\mathbb{B}_{1/3}(0)) = \mathbb{R} \smallsetminus \mathbb{Q}$. *Thus* $\chi_{\mathbb{Q}}$ *cannot be continuous at any point of* $\mathbb{R}$ *because neither* $\mathbb{Q}$ *nor* $\mathbb{R} \smallsetminus \mathbb{Q}$ *contains an interval. (**Warning!!** Since the concept of continuity is strongly linked to openness and closure (by definition), you need to be extremely careful when discussing continuity of a function on a given space. For instance, in this case we have that* $\chi_{\mathbb{Q}}$ *has no points of continuity relative to* $\mathbb{R}$, *but the restriction of* $\chi_{\mathbb{Q}}$ *to* $\mathbb{Q}$ *is everywhere continuous relative to* $\mathbb{Q}$ *itself! This is because for any given space M, the whole space M is always open. It may be the case that M is contained in an ambient space* $\widetilde{M}$ *and that M is closed relative to* $\widetilde{M}$, *but M will always be open (and closed!) relative to itself.)*

*b) A function* $f \colon M \to N$ *between metric spaces M and N is called an **isometry** if f preserves distances, that is, if* $\rho(f(x), f(y)) = d(x, y)$ *for all* $x, y \in M$. *Obviously, an isometry is continuous. The natural inclusions from* $\mathbb{R}$ *into* $\mathbb{R}^2$ *(i.e.,* $x \mapsto (x, 0)$*) and from* $\mathbb{R}^2$ *into* $\mathbb{R}^3$ *(i.e.,* $(x, y) \mapsto (x, y, 0)$*) are examples of injective isometries.*

*c) Let* $f \colon \mathbb{N} \to \mathbb{R}$ *be any function. Then f is always continuous! Why? Because* $\{n\}$ *is an open ball in* $\mathbb{N}$. *Specifically,* $\{n\} = \mathbb{B}_{1/2}(n) \subset f^{-1}(\mathbb{B}_\varepsilon(f(n)))$ *for any* $\varepsilon > 0$. *(This example is topological in nature: all maps from spaces that are endowed with a discrete topology are always continuous, because they have "too many open sets." This is what in topology is called a **fine topology**. Discrete topologies are the finest of them all, i.e., they have more open sets than any other possible topology. These words may not mean much to you now, but it will all start to make sense when we get to cover the beautiful subject of topology on Chapter 3.)*

*d) $f\colon \mathbb{R} \to \mathbb{N}$ is continuous if and only if $f$ is constant! Why? [Hint: Recall that $\mathbb{R}$ has no nontrivial clopen sets.] Note that this case is the complete opposite of the situation described in c): While every map <u>from</u> a discrete space is always continuous, very few maps (i.e., only the constant-valued ones) <u>to</u> a discrete space are continuous.*

*e) If $y$ is any fixed element of $(M, d)$, then the real-valued function $f(x) = d(x, y)$ is continuous on M.*

**Definition 25.** *Given a nonempty set $A$ and a point $x \in M$, we define the **distance** from $x$ to $A$ by $d(x, A) = \inf\{d(x, a) \mid a \in A\}$.*

Clearly, $0 \le d(x, A) < \infty$ for any $x$ and any $A$, but it is not necessarily true that $d(x, A) > 0$ when $x \notin A$. For instance, $d(x, \mathbb{Q}) = 0$ for any $x \in \mathbb{R}$.

**Proposition 14.** *The distance from any point $x$ to a set $A$ is $0$, i.e., $d(x, A) = 0$, if and only if $x \in \bar{A}$.*

*Proof.* The proof is quite short. Simply notice that $d(x, A) = 0$ if and only if there is a sequence of points $\{a_n\}$ in $A$ such that $d(x, a_n) \to 0$. But this in turn means that $a_n \to x$ and, thus $x \in \bar{A}$, as desired. $\qquad\square$

Note that this proposition has given us another connection between limits in $M$ and limits in $\mathbb{R}$. Loosely speaking, the statement shows that $0$ is a limit point of $\{d(x, a) \mid a \in A\}$ if and only if $x$ is a limit point of $A$. We can even squeeze more juice out of this observation by checking that the map $x \mapsto d(x, A)$ is actually continuous. For this, it suffices to establish the inequality of the following proposition.

**Proposition 15.** *For any points $x, y$ and any set $A$, we have that $|d(x, A) - d(y, A)| \le d(x, y)$.*

*Proof.* It is true by the triangle inequality that $d(x, a) \le d(x, y) + d(y, a)$ for any $a \in A$. But $d(x, A)$ is a lower bound for $d(x, a)$; hence $d(x, A) \le d(x, y) + d(y, a)$. Now, by taking the infimum over $a \in A$, we get $d(x, A) \le d(x, y) + d(y, A)$. Since $d$ is symmetric, the roles of $x$ and $y$ are interchangeable and so we are done. $\qquad\square$

To appreciate what this has done for us, let's make two simple observations. First, if $f\colon M \to \mathbb{R}$ is a continuous function, then the set $E = \{x \in M \mid f(x) = 0\}$ is closed, since $\{0\}$ is closed in $\mathbb{R}$ and so its inverse image $E$ must be closed as well by continuity of $f$. Conversely, if $E$ is a closed set in $M$, then $E$ is the **zero set** of some continuous real-valued function on $M$; in particular, $E = \{x \in M \mid d(x, E) = 0\}$. Thus a set $E$ is closed if and only if $E = f^{-1}(\{0\})$ for some continuous function $f\colon M \to \mathbb{R}$. The moral of this story is that if you know all of the open (and hence closed) sets in a metric space $M$, then you know all of the continuous real-valued functions on $M$, and vice versa.

---

## 1.8 Connectedness, Completeness, & Compactness

### 1.8.1 Connected Sets

Recall from your glory days of kindergarten calculus the well known *Intermediate Value Theorem*, which is the formal statement of the informal notion that the graph of a continuous function is "unbroken." In this subsection we are going to analyze a generalization of this result that extends to metric spaces (in fact, it also applies to the even more general case of *topological spaces*, which you will study on Chapter 3).

**Definition 26.** *If A is a subspace of an ambient space M so that $A \subsetneq M$ (this notation means that A is contained in M, but it is not equal to the entire space M),[7] then we say that A is a **proper subspace** of M. If, moreover, A is nonempty, then we say that A is a **nontrivial subspace** of M.[8]*

**Definition 27.** *We say that a metric space M is **disconnected** if M can be split into the disjoint union of two nontrivial open sets, that is, if there are nonempty open sets U and V in M with $U \cap V = \varnothing$ and $U \cup V = M$ (note that such a disjoint union is usually denoted in the literature*

---

[7] Normally, the notation "$A \subset M$" suffices, since it's technically the same as "$A \subsetneq M$," but nonetheless you may encounter the latter quite often when we want to emphasize that the subset/subspace is properly contained.

[8] Be aware that some authors require that a proper subspace be nonempty by definition.

*as $U \amalg V$ or $U \sqcup V$). The open sets $U$ and $V$ are called the **connected components** (or just the **components**) of M, and the pair $(U, V)$ is referred to as the **disconnection** of M; we say that M is **connected** (duh!) if no such disconnection can be found. Note also that a disconnection doesn't have to be just a pair; a space can be composed of the disjoint union of many (in fact, even uncountably many!) components (see Figure 1.4 for an example of a space N with three components).*
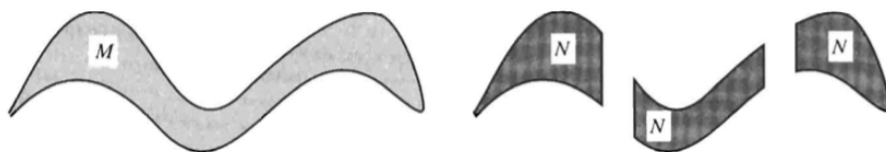


Figure 1.4: Illustration of a connected set (M) and a disconnected one (N).

**Remark:** Note that Definition 27 could be just as easily formulated in terms of closed sets instead. After all, if a disconnection $(U, V)$ exists, then the disconnecting sets $U$ and $V$ are also closed: $U = V^c$ and $V = U^c$. In other words, $U$ and $V$ are clopen sets. Conversely, if M contains a nontrivial clopen subset $U$ (i.e. a clopen set other than $\varnothing$ or M), then $(U, U^c)$ is a disconnection for M. This gives us another way to present connectedness; in fact, it is this formulation the one that we often use in proofs: M is connected if and only if it contains no nontrivial clopen sets.

**Example 20.** *a) Any interval $[a, b]$ is connected; in fact the entire real line $\mathbb{R}$ is connected. (We can show that $\mathbb{R}$ contains no nontrivial clopen sets by showing that if $U$ is a nontrivial open subset of $\mathbb{R}$, then $\overline{U} \supsetneq U$.)*

*b) A discrete space containing two or more points is disconnected.*

*c) Both the empty set $\varnothing$ and any one-point space are always connected (by default).*

*d) The Cantor set $\Delta$ is (very!) disconnected. Indeed, for any with $x < y$ there is a $z \notin \Delta$ such that $x < z < y$. Thus, $\Delta$ is disconnected by the (relatively) open sets $U = [0, z) \cap \Delta$ and $V = (z, 1] \cap \Delta$. The Cantor set is an example of a **totally disconnected space**; in the more general setting of topology, which we will discuss in Chapter 3, a topological space X is said to be **totally disconnected** if the connected components in X are the one-point sets.* 🌍

**Theorem 42.** *If M is connected and the mapping $\varphi\colon M \to N$ is continuous and onto, then N is connected as well. In particular (even when $\varphi$ is not onto), the continuous image of a connected space is connected.*

*Proof.* If $U$ is a clopen nontrivial subset of $N$, according to the open and closed set conditions for continuity, the inverse image $\varphi^{-1}(U)$ must also be clopen in $M$. Since $\varphi$ is onto and $U \neq \varnothing$, it follows that $\varphi^{-1}(U)$ is also nonempty. Similarly, $\varphi^{-1}(U^c) \neq \varnothing$. Therefore $\varphi^{-1}(U)$ is a nontrivial clopen subset of $M$, contrary to our assumption that $M$ is connected. Thus we have reached a contradiction, from which follows that the existence of such a clopen set $U$ in $N$ is impossible; hence $N$ is connected just like $M$. $\qquad\square$

**Corollary 9.** *If M is homeomorphic to N and M is connected, then N is also connected.*[9]

*Proof.* Since $N$ is the continuous image of $M$ under the isomorphism, the result follows immediately from Theorem 42. $\qquad\square$

**Corollary 10** (**Generalized Intermediate Value Theorem**). *Let M be connected. Then every continuous real-valued function $f\colon M \to \mathbb{R}$ has the intermediate value property.*

*Proof.* If $f$ assumes values $\alpha, \beta \in \mathbb{R}$ such that $\alpha < \beta$, but it fails to assume some value $\xi$ between $\alpha$ and $\beta$ (i.e., $\alpha < \xi < \beta$), then we have that

$$U = \{x \in M \mid f(x) < \xi\} \qquad \text{and} \qquad V = \{x \in M \mid f(x) > \xi\}$$

separate $M$, i.e., $M = U \amalg V$. But this contradicts our assumption that $M$ is connected. $\quad\square$

**Theorem 43.** *The union of connected sets sharing a common point is connected.*

*Proof.* Let $S = \cup S_\alpha$, where each $S_\alpha$ is connected and they have a common pint $p \in \cap S_\alpha$. If the union $S$ were disconnected, it would have a separation $(U, U^c)$, where $U$ and $U^c$ are both nontrivial and clopen. One of them contains $p$; say, WLOG, that it is $U$. Then $U \cap S_\alpha$ is a nonempty clopen subset of $S_\alpha$. But since $S_\alpha$ is connected by assumption, we then must

---

[9]  As you will see in Chapter 3, connectedness is a "topological property," i.e., it is a property that is "invariant" (i.e., preserved) by homeomorphisms.

have $U \cap S_\alpha = S_\alpha$ for each $\alpha$, and $U = S$. This in turn implies that $U^c = \varnothing$, contrary to our assumption that $U^c$ is nontrivial.                                                                $\square$

**Definition 28.** *Given a metric space M and points $p, q \in M$, a **path** joining p to q is a continuous map $\gamma \colon [a,b] \subset \mathbb{R} \to M$ such that $\gamma(a) = p$ and $\gamma(b) = q$. (Note that we may also consider the **reversed path** $\overline{\gamma} \colon [a,b] \subset \mathbb{R} \to M$ that satisfies $\overline{\gamma}(a) = q$ and $\overline{\gamma}(b) = p$; in other words, $\overline{\gamma}$ has the same image in M as $\gamma$, but it has the reversed orientation.) If each pair of points in M can be joined by a path in M, then we say that M is **path-connected** (see Figure 1.5).*
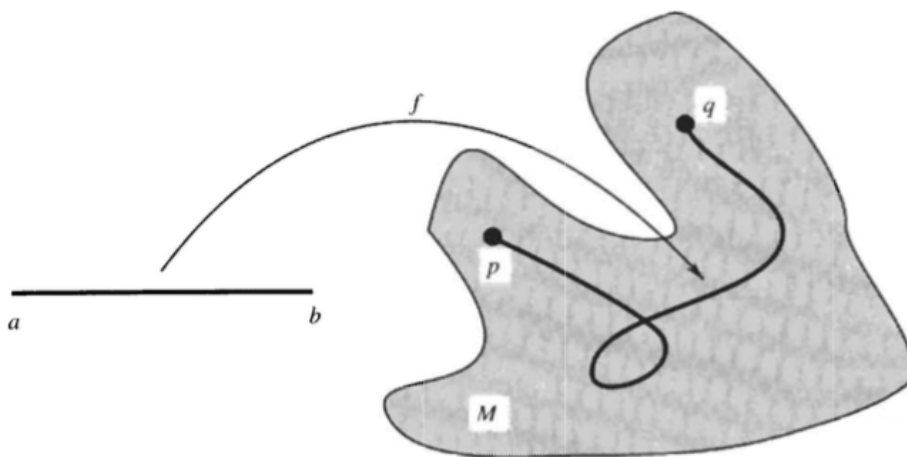


Figure 1.5: A path $f \colon [a,b] \to M$ joining $p$ and $q$.

**Theorem 44.** *Path connectedness implies connectedness.*

*Proof.* Let $M$ be a path-connected space and assume, to the contrary, that $M$ is not connected. Then $M = U \amalg U^c$ for some clopen nontrivial subset $U \subsetneq M$. Now choose $p \in U$ and $q \in U^c$; by the assumption that $M$ is path-connected, here exists a path $\gamma \colon [a,b] \subset \mathbb{R} \to M$ joining $p$ to $q$. But then the inverse images $\gamma^{-1}(U)$ and $\gamma^{-1}(U^c)$ contradict the connectedness of the interval $[a,b]$.                                                                $\square$

Note however, that the converse of Theorem 44 does not hold in general:

**Example 21 (The Topologist's Sine Curve).** *The **topologist's sine curve** (shown in Figure 1.6), is a compact (a concept to be described later in this section) and connected set that is **not** path-connected. It is a metric space $(M, d)$ with d being the usual Euclidean distance and M given by the union $M = \Gamma \amalg Y$, where*

$$\Gamma = \left\{ (x, y) \in \mathbb{R}^2 \mid y = \sin 1/x, \, 0 < x \le 1/\pi \right\},$$
$$Y = \left\{ (0, y) \in \mathbb{R}^2 \mid -1 < yx \le 1 \right\}.$$

*Take another look at Figure 1.6. Does M seem connected to you? You bet your sorry a\*\* it is! How*
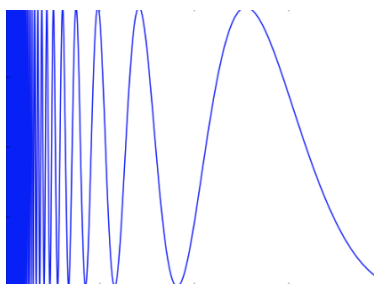


Figure 1.6: The topologist's sine curve.

*do we prove this? Note that the graph $\Gamma$ is connected and M is its closure. Now, since the closure of a connected space is connected[10] (this is not so trivial; prove it!), it follows that M must be connected, as desired.*

### 1.8.2 Totally Bounded Sets

**Definition 29.** *A subset A in a metric space $(M, d)$ is said to be **totally bounded** if, given any $\varepsilon > 0$, there exist finitely many points $x_1, \ldots, x_n \in M$ such that $A \subset \cup_{i=1}^{n} \mathbb{B}_\varepsilon(x_i)$ (when A is contained in such union, we often say that A is **covered** by finitely many $\varepsilon$-balls). Loosely speaking, each point $x \in A$ is within a distance of $\varepsilon$ from some other point $x_i$. For this reason, we often refer to the set $\{x_1, \ldots, x_n\}$ as being $\varepsilon$-**dense** in A (you may also find in the literature that $\{x_1, \ldots, x_n\}$ is called an $\varepsilon$-**net** for A).*

---

[10] In fact, more generally, if $S$ and $\widetilde{S}$ are any two sets such that $S \subset \widetilde{S} \subset \overline{S}$, and $S$ is connected, then so is $\widetilde{S}$.

**Lemma 8.** *A subset $A$ in a metric space $(M, d)$ is totally bounded if and only if, given $\varepsilon > 0$, there are finitely many sets $A_1, \ldots, A_n \subset A$, with $\operatorname{diam} A_i < \varepsilon$ for all $i$, such that $A \subset \cup_{i=1}^n A_i$.*

*Proof.* $(\Rightarrow)$ First we suppose that $A$ is totally bounded. Given $\varepsilon > 0$, we may choose $x_1, \ldots, x_n \in M$ such that $A \subset \cup_{i=1}^n \mathbb{B}_\varepsilon(x_i)$. Then $A$ is covered by the sets $A_i = A \cap \mathbb{B}_\varepsilon(x_i) \subset A$ and $\operatorname{diam} A_i \leq \operatorname{diam} \mathbb{B}_\varepsilon(x_i) \leq 2\varepsilon$ for each $i$.

$(\Leftarrow)$ Conversely, given $\varepsilon > 0$, suppose that there are finitely many sets $A_1, \ldots, A_n \subset A$, with $\operatorname{diam} A_i < \varepsilon$ for all $i$, such that $A \subset \cup_{i=1}^n A_i$. Given $x_i \in A_i$, we have $A_i \subset \mathbb{B}_{2\varepsilon}(x_i)$ for each $i$ and thus, $A \subset \cup_{i=1}^n A_i \subset \cup_{i=1}^n \mathbb{B}_{2\varepsilon}(x_i)$ . Since $\varepsilon$ is arbitrary in either case, we are done. $\qquad\square$

**Example 22.** *a)* By the triangle inequality, a totally bounded set is necessarily bounded (make sure you understand why this is true). Note also that any subset of a totally bounded set is again totally bounded.

*b)* A finite set is always totally bounded. In a discrete space, a set is totally bounded if and only if it is finite. Thus, for instance $\mathbb{N}$ and $\mathbb{Z}$ are not totally bounded, but $\mathbb{Z}_n$[11] (for any $n \in \mathbb{N}$) is totally bounded.

*c)* Any given subset of $\mathbb{R}$ is totally bounded if and only if it is bounded: in particular, if $I$ is a closed and bounded interval in $\mathbb{R}$ and $\varepsilon > 0$, then $I$ can be covered by finitely many closed subintervals $J_1, \ldots, J_n$, each of length at most $\varepsilon$ (prove this!). Thus we have that total boundedness is not a topological property (in fact, total boundedness depends intimately on the metric at hand): as we briefly discussed on Subsection *1.7.1*, topological properties are those properties of a space that are preserved by bicontinuous mappings called "homeomorphisms." Now, $\mathbb{R}$ is homeomorphic to the interval $(0, 1)$, and the former is not bounded whereas the latter is. Thus, since a subset of $\mathbb{R}$ is totally bounded if and only if it is bounded, it follows that total boundedness is not preserved by homeomorphisms and thus it is not a topological property.

*d)* This feature of $\mathbb{R}$ that we presented on part c) does not in general hold for arbitrary spaces. We will now show that not every bounded set is totally bounded; the discrete metric gives us a clue as to how we might construct such a set: Recall the sequence $\{e^{(n)}\} \in \ell_1$ in which $e^{(j)} = (0, \ldots, 0, 1, 0, \ldots)$, where the single nonzero entry is in the $j^{th}$ place. Then, $\{e^{(n)}\}$ is a bounded set in $\ell_1$ (since $\|e^{(j)}\|_1 = 1$ for all $j$), but it is not totally bounded. Why? Because $\|e^{(j)} - e^{(i)}\|_1 = 2$ for

---

[11] This is the set of all integers modulus $n$; if you're not already familiar with these algebraic constructions don't worry, you'll learn about them in the next chapter.

$j \neq i$; thus, $\{e^{(n)}\}$ cannot be covered by finitely many balls of radius $< 2$. In fact, $\{e^{(n)}\}$ is discrete in its relative (i.e., subspace) metric.

Now we give a sequential criterion for total boundedness:

**Lemma 9.** *Let $\{x_n\}$ be a sequence in $(M, d)$, and let $A \subset M$ be the range of this sequence.*

    *a) If $\{x_n\}$ is Cauchy, then $A$ is totally bounded.*

    *b) If $A$ is totally bounded, then $\{x_n\}$ has a Cauchy subsequence.*

*Proof of a).* Let $\varepsilon > 0$. Then, since $\{x_n\}$ is assumed to be Cauchy, there is some index $N \geq 1$ such that diam $\{x_n \mid n \geq N\} < \varepsilon$. Thus,

$$A = \underbrace{\{x_1\} \cup \cdots \cup \{x_{N-1}\} \cup \{x_n \mid n \geq N\}}_{N \text{ sets of diameter } < \varepsilon}. \qquad \square$$

*Proof of b).* If we let $A$ be finite, the result is trivial. Thus, assume that $A$ is an infinite totally bounded set. Then $A$ can be covered by finitely many sets of diameter $< 1$, where at least one of these sets must contain infinitely many points of $A$. Call this set $A_1$. But then $A_1$ is also totally bounded, and so it can be covered by finitely many sets of diameter $< 1/2$. One of these, call it $A_2$, contains infinitely many points of $A_1 \ldots$

Continuing this process, we find a decreasing sequence of sets $A \supset A_1 \supset A_2 \supset \ldots$, where each $A_k$ contains infinitely many $x_n$ and where diam $A_k < 1/k$. In particular, we may choose a subsequence $\{x_{n_k}\}_k$ with $x_{n_k} \in A_k$ for all $k$. The fact that $\{x_{n_k}\}_k$ is Cauchy is now clear since diam $\{x_{n_j} \mid j \geq k\} \leq$ diam $A_k < 1/k$. $\qquad \square$

**Example 23.** *a) The sequence given by $x_n = (-1)^n$ in $\mathbb{R}$ shows that a Cauchy subsequence is the best that we can hope for in part ii) of Lemma 9.*

    *b) Note that the sequence $\{e^{(n)}\} \in \ell_1$ that we have previously discussed has no Cauchy subsequence.*

Now we are truly ready for our sequential characterization of total boundedness:

**Theorem 45.** *A set A is totally bounded if and only if every sequence in A has a Cauchy subsequence.*

*Proof.* ($\Rightarrow$) This is crystal clear from Lemma 9.

($\Leftarrow$) Suppose, to the contrary, that $A$ is not totally bounded. Then, there is some $\varepsilon > 0$ such that $A$ cannot be covered by finitely many $\varepsilon$-balls. Thus, by induction, we can find a sequence $\{x_n\}$ in $A$ such that $d(x_n, x_m) \geq \varepsilon$ whenever $m \neq n$. But then, $\{x_n\}$ has no Cauchy subsequence, which contradicts our assumption. $\square$

**Corollary 11 (The Bolzano-Weierstrass Theorem).** [12] *Every bounded infinite subset of $\mathbb{R}$ has a limit point in $\mathbb{R}$.*

*Proof.* Let $A$ be a bounded infinite subset of $\mathbb{R}$. Then, in particular, there is a sequence $\{x_n\}$ of distinct points in $A$. Since $A$ is totally bounded, there is a Cauchy subsequence $\{x_{n_k}\}_k$ of $\{x_n\}$. But Cauchy sequences in $\mathbb{R}$ converge, and so $\{x_{n_k}\}_k$ converges to some $x \in \mathbb{R}$. Thus, $x$ is a limit point of $A$. $\square$

### 1.8.3  Complete & Compact Metric Spaces

We are almost ready to discuss *compact* sets, but first we need a few results about *completeness*.

**Definition 30.** *A metric space M is said to be **complete** if every Cauchy sequence in M converges to a point in M.*

**Example 24.** *a)* $\mathbb{R}$ *is complete. This is a consequence of the least upper bound axiom; in fact, as we will see, the completeness of $\mathbb{R}$ is actually equivalent to the least upper bound axiom.*

*b)* $\mathbb{R}^n$ *is complete (because $\mathbb{R}$ is).*

*c) Any discrete space is complete (trivially).*

---

[12] There are other incarnations of this theorem; for instance, another formulation says that every bounded sequence of real numbers has a convergent subsequence.

*d) $(0,1)$ is not complete; hence, completeness is not preserved by homeomorphisms. (recall that $(0,1) \cong \mathbb{R}$)*

*e) $\ell_1$, $\ell_2$, $\ell_p$ and $\ell_\infty$ are all complete. We will sketch the proof for $\ell_2$ below (the rest of the proofs are all very similar).*

The proof that $\ell_2$ is complete is based on a few simple principles that will generalize to all sorts of different settings. This generality will become all the more apparent if we introduce a slight change in our notation. Since a *sequence* is just another name for a function on $\mathbb{N}$, let's agree to write an element $f \in \ell_2$ as $f = \{f(k)\}$, in which case

$$\|f\|_2 = \left( \sum_{k=1}^{\infty} |f(k)|^2 \right)^{1/2}.$$

For example, the vectors $e^{(n)}$ will now be written as $e_n$, where $e_n(k) = \delta_{n,k}$. In case you're not familiar with $\delta_{n,k}$, it is known as **Kronecker's delta**, which is defined by

$$\delta_{n,k} = \begin{cases} 1 & \text{if } n = k, \\ 0 & \text{otherwise.} \end{cases}$$

Now let $\{f_n\}$ be a sequence in $\ell_2$, where now we write $f_n = \{f_n(k)\}_k$, and suppose that $\{f_n\}$ is Cauchy in $\ell_2$. That is, suppose that for each $\varepsilon > 0$ there exists an $n_0$ such that $\|f_n - f_m\|_2 < \varepsilon$ whenever $m, n \geq n_0$. Of course, we want to show that $\{f_n\}$ converges, in the metric of $\ell_2$, to some $f \in \ell_2$. We will break the proof into three steps:

**Step 1** $f(k) = \lim_{n \to \infty} f_n(k)$ exists in $\mathbb{R}$ for each $k$.

To see why, note that $|f_n(k) - f_m(k)| \leq \|f_n - f_m\|_2$ for any $k$, and hence $\{f_n(k)\}_k$ is Cauchy in $\mathbb{R}$ for each $k$. Thus, $f$ is the obvious candidate for the limit of $\{f_n\}$, but we still have to show that the convergence takes place in the metric space $\ell_2$; that is, we need to show that $f \in \ell_2$ and that $\|f_n - f\|_2 \to 0$ (as $n \to \infty$).

**Step 2** $f \in \ell_2$; that is, $\|f\|_2 < \infty$.

We know that $\{f_n\}$ is bounded in $\ell_2$ (why?); say $\|f_n\|_2 \leq B$ for all $n$. Thus, for any fixed $N < \infty$, we have:

$$\sum_{k=1}^{N} |f(k)|^2 = \lim_{n \to \infty} \sum_{k=1}^{N} |f_n(k)|^2 \leq B^2.$$

Since this holds for any $N$, we get that $\|f\|_2 \leq B$.

Step 3 Now we repeat Step 2 (more or less) to show that $f_n \to f$ in $\ell_2$.

Given $\varepsilon > 0$, choose $n_0$ so that $\|f_n - f_m\|_2 < \varepsilon$ whenever $m, n \geq n_0$. Then, for any $N$ and any $n \geq n_0$,

$$\sum_{k=1}^{N} |f(k) - f_n(k)|^2 = \lim_{m\to\infty} \sum_{k=1}^{N} |f_m(k) - f_n(k)|^2 \leq \varepsilon^2.$$

Since this holds for any $N$, we have $\|f - f_n\|_2 \leq \varepsilon$ for all $n \geq n_0$. That is, $f_n \to f$ in $\ell_2$, as desired.                                                                      ∎

**Example 25.** *a) Just having a candidate for a limit is not enough. Consider the sequence $\{f_n\}$ in $\ell_\infty$ defined by $f_n = (1, \ldots, 1, 0, \ldots)$, where the first $n$ entries are 1 and the rest are 0. The "obvious" limit is $f = (1, 1, \ldots)$ (all 1's), but $\|f - f_n\|_\infty = 1$ for all $n$ (what's wrong?).*

*b) Worse still, sometimes the "obvious" limit is not even in the space. Consider the same sequence as in part a) and note that each $f_n$ is actually an element of $c_0$. This time, the natural candidate $f$ is not in $c_0$ (again, what's wrong?).*                                                                      ☙

As you can see, there can be a lot of details to check in a proof of completeness, and it would be handy to have at least a few easy cases available. For example: when is a subset of a complete space complete? The answer is given in the following theorem.

**Theorem 46.** *Let $(M, d)$ be a complete metric space and let $A$ be a subset of M. Then, $(A, d)$ is complete if and only if A is closed in M.*

*Proof.* $(\Rightarrow)$ First suppose that $(A, d)$ is complete, and let $\{x_n\}$ be a sequence in $A$ that converges to some point $x \in M$. Then $\{x_n\}$ is Cauchy in $(A, d)$ and so it converges to some point of $A$. That is, we must have $x \in A$ and, hence, $A$ is closed.

$(\Leftarrow)$ Next suppose that $\{x_n\}$ is a Cauchy sequence in $(A, d)$. Then $\{x_n\}$ is also Cauchy in $(M, d)$. Hence, we have that $\{x_n\}$ converges to some point $x \in M$. But $A$ is closed and so, in fact, $x \in A$. Thus, $(A, d)$ is complete.                                                                      □

**Example 26.** *a) $[0, 1]$, $[0, \infty)$, $\mathbb{N}$, and the Cantor set $\Delta$ are all complete.*

*b) It follows from Theorem 46 that if a metric space $(M, d)$ is both complete and totally bounded, then every sequence in M has a convergent subsequence. In particular, any closed, bounded subset of*

$\mathbb{R}$ *is both complete and totally bounded. Thus, for example, every sequence in* $[a, b]$ *has a convergent subsequence. As you can easily imagine, the interval* $[a, b]$ *is a great place to do analysis! We will pursue the consequences of this felicitous combination of properties when we explore compact sets.* ✪

Our next result underlines the fact that complete spaces have a lot in common with $\mathbb{R}$.

**Theorem 47.** *For any metric space* $(M, d)$, *the following statements are equivalent:*

  i)  $(M, d)$ *is complete.*

 ii)  (THE NESTED SET THEOREM) *Let* $F_1 \supset F_2 \supset \ldots$ *be a decreasing sequence of nonempty closed sets in* $M$ *with* $\mathrm{diam}(F_n) \to 0$. *Then,* $\bigcap_{n=1}^{\infty} F_n \neq \varnothing$ *(in fact, it contains exactly one point).*

iii)  (THE BOLZANO-WEIERSTRASS THEOREM) *Every infinite, totally bounded subset of* $M$ *has a limit point in* $M$.

*Proof.*  ($i$) $\implies$ $ii$)) Given $\{F_n\}$ as in ii), choose $x_n \in F_n$ for each $n$. Then, since the $F_n$ decrease, $\{x_k \mid k \geq n\} \subset F_n$ for each $n$, and hence $\mathrm{diam}\{x_k \mid k \geq n\} \to 0$ as $n \to \infty$. That is, $\{x_n\}$ is Cauchy. Since $M$ is complete, we have $x_n \to x$ for some $x \in M$. But the $F_n$ are closed, and so we must have $x \in F_n$ for all $n$. Thus, $\bigcap_{n=1}^{\infty} F_n \neq \varnothing$.

($ii$) $\implies$ $iii$)) Let $A$ be an infinite, totally bounded subset of $M$. Recall that we have shown that $A$ contains a Cauchy sequence $\{x_n\}$ comprised of distinct points ($x_n \neq x_m$ for $n \neq m$). Now, setting $A_n = \{x_k \mid k \geq n\}$, we get $A \supset A_1 \supset A_2 \supset \ldots$, each $A_n$ is nonempty (even infinite), and $\mathrm{diam}(A_n) \to 0$. That is, ii) *almost* applies, but clearly, $\overline{A}_n \supset \overline{A}_{n+1} \neq \varnothing$ for each $n$, and $\mathrm{diam}(\overline{A}_n) = \mathrm{diam}(A_n) \to 0$ as $n \to \infty$. Thus there exists an $x \in \bigcap_{n=1}^{\infty} \overline{A}_n \neq \varnothing$. Now $x_n \in A_n$ implies that $d(x_n, x) \leq \mathrm{diam}(\overline{A}_n) \to 0$. That is, $x_n \to x$ and thus $x$ is a limit point of $A$.

($iii$) $\implies$ $i$)) Let $\{x_n\}$ be Cauchy in $(M, d)$. We just need to show that $\{x_n\}$ has a convergent subsequence. Now, the set $A = \{x_n \mid n \geq 1\}$ is totally bounded (prove it!). If $A$ happens to be finite, we are done (why?). Otherwise, iii) tells us that $A$ has a limit point $x \in M$. It follows that some subsequence of $\{x_n\}$ converges to $x$. □

We are finally ready to study compact metric spaces.

**Definition 31.** *A metric space* $(M, d)$ *is said to be* **compact** *if it is both complete and totally bounded.*

**Example 27.** *a) A subset K of* $\mathbb{R}$ *is compact if and only if K is closed and bounded. This fact is usually referred to as the* **Heine-Borel theorem**. *Hence, a closed bounded interval* $[a, b]$ *is compact. Also, the Cantor set* $\Delta$ *is compact. The interval* $(0, 1)$, *on the other hand, is not compact.*

*b) A subset K of* $\mathbb{R}^n$ *is compact if and only if K is closed and bounded.*

*c) It is important that we not confuse the first two examples with the general case. Recall that the set* $\{e_n \mid n \geq 1\}$ *is closed and bounded in* $\ell_\infty$ *but not totally bounded –hence not compact. Taking this a step further, notice that the closed ball* $\{x \mid \|x\|_\infty \leq 1\}$ *in* $\ell_\infty$ *is not compact, whereas any closed ball in* $\mathbb{R}^n$ *is compact.*

*d) A subset of a discrete space is compact if and only if it is finite. (why?)*

Just as with completeness and total boundedness, we will want to give several equivalent characterizations of compactness. In particular, since neither completeness nor total boundedness is preserved by homeomorphisms, our newest definition does not appear to be describing a topological property. Let's remedy this immediately by giving a sequential characterization of compactness that will turn out to be invariant under homeomorphisms:

**Theorem 48.** $(M, d)$ *is* **(sequentially) compact** *if and only if every sequence in M has a subsequence that converges to a point in M.*

*Proof.* The proof follows immediately from the fact that a space $M$ is totally bounded and complete (hence compact) if and only if every sequence in $M$ has a Cauchy subsequence (that converges to a point in $M$). $\square$

Compactness is quite a valuable property to have available on a space, as convergent sequences are easy to come by in a compact space. In the case that you are dealing with a non-convergent sequence, simply extract a subsequence that does converge and use that one instead. That's the best that one could hope for! Given a compact space, it is easy to decide which of its subsets are compact: Every bounded infinite subset of $\mathbb{R}$ has a limit point in $\mathbb{R}$.

**Corollary 12.** *Let A be a subset of a metric space M. If A is compact, then A is closed in M. If M is compact and A is closed, then A is compact.*

*Proof.* Suppose that $A$ is compact, and let $\{x_n\}$ be a sequence in $A$ that converges to a point $x \in M$. Then, from the above theorem, $\{x_n\}$ has a subsequence that converges in $A$, and hence we must have $x \in A$. Thus, $A$ is closed.

Next, suppose that $M$ is compact and that $A$ is closed in $M$. Given an arbitrary sequence $\{x_n\}$ in $A$, Theorem 48 supplies a subsequence of $\{x_n\}$ that converges to a point $x \in M$. But since $A$ is closed, we must have $x \in A$. Thus, $A$ is compact. □

To show that compactness is indeed a topological property, let us show that the continuous image of a compact set is again compact:

**Theorem 49.** *Let $f : (M, d) \to (N, \rho)$ be continuous. If K is compact in M, then $f(K)$ is compact in N.*

*Proof.* Let $\{y_n\}$ be a sequence in $f(K)$. Then, $y_n = f(x_n)$ for some sequence $\{x_n\}$ in $K$. But, since $K$ is compact, $\{x_n\}$ has a convergent subsequence, say, $x_{n_k} \to x \in K$. Then, since $f$ is continuous, $y_{n_k} = f(x_{n_k}) \to f(x) \in f(K)$. Thus, $f(K)$ is compact. □

The theorem above gives us a wealth of useful information. In particular, it tells us that real-valued continuous functions on compact spaces are quite well behaved:

**Corollary 13.** *Let $(M, d)$ be compact. If $f : M \to \mathbb{R}$ is continuous, then $f$ is bounded. Moreover, $f$ attains its maximum and minimum values.*

*Proof.* The image $f(M)$ is compact in $\mathbb{R}$; hence it is closed and bounded. Moreover, $\sup f(M)$ and $\inf f(M)$ are actually elements of $f(M)$. That is, there exist $x, y \in M$ such that $f(x) \leq f(t) \leq f(y)$ for all $t \in M$. (In this case we would write $f(x) = \min_{t \in M} f(t)$ and $f(y) = \max_{t \in M} f(t)$). □

**Corollary 14.** *If $f : [a, b] \to \mathbb{R}$ is continuous, then the range of $f$ is a compact interval $[c, d]$ for some $c, d \in \mathbb{R}$.*

**Corollary 15.** *If M is a compact metric space, then $\|f\|_\infty = \max_{t \in M} |f(t)|$ defines a norm on $C^0(M)$, the vector space of continuous real-valued functions on M.*

It appears that compactness is the analogue of "finite." To better appreciate this, we will need a slightly more esoteric characterization of compactness. A bit of preliminary detail-checking will ease the transition.

**Lemma 10.** *In a metric space M, the following are equivalent:*

 *i) If $\mathcal{G}$ is any collection of open sets in M with $\bigcup\{G \mid G \in \mathcal{G}\} \supset M$, then there are finitely many sets $G_1, \ldots, G_n \in \mathcal{G}$ with $\bigcup_{i=1}^{n} G_i \supset M$.*

 *ii) If $\mathcal{F}$ is any collection of closed sets in M such that $\bigcap_{i=1}^{n} F_i \neq \emptyset$ for all choices of finitely many sets $F_1, \ldots, F_n \in \mathcal{F}$, then $\bigcap\{F \mid F \in \mathcal{F}\} \neq \emptyset$.*

The first condition is the precise mathematical expression of the phrase ***every open cover has a finite subcover***, which is often used as the definition of compactness. The second condition is abbreviated by saying ***every collection of closed sets with the finite intersection property has nonempty intersection***. These may at first seem to be unwieldy statements to work with, but each is worth the trouble. Here is why we care:

- Condition i) implies that $M$ is totally bounded because, for any $\varepsilon > 0$, the collection $\mathcal{G} = \{B_\varepsilon(x) \mid x \in M\}$ is an open cover for $M$.

- Condition ii) implies that $M$ is complete because it easily implies the nested set theorem (if $F_1 \supset F_2 \supset \ldots$ are nonempty, then $\bigcap_{i=1}^{n} F_i = F_n \neq \emptyset$).

Putting these two conditions together we have our new characterization of compactness:

**Theorem 50.** *M is compact if and only if it satisfies either (hence both, since they are equivalent) conditions i) and ii) in Lemma 10.*

*Proof.* ($\Leftarrow$) As noted above, conditions i) and ii) imply that $M$ is totally bounded and complete, hence compact.

($\Rightarrow$) We need to show that compactness will imply, say, i) (there is no loss of generality here, since both conditions i) and i)) are equivalent). To this end, suppose that $M$ is compact,

and suppose that $\mathcal{G}$ is an open cover for $M$ that admits no finite subcover. We will work toward a contradiction:

$M$ is totally bounded, so $M$ can be covered by finitely many closed sets of diameter at most 1. It follows that at least one of these, call it $A_1$, cannot be covered by finitely many sets from $\mathcal{G}$. Certainly $A_1 \neq \emptyset$ (since the empty set is easy to cover!). Note that $A_1$ must be infinite.

Next, $A_1$ is totally bounded, so $A_1$ can be covered by finitely many closed sets of diameter at most $1/2$. At least one of these, call it $A_2$, cannot be covered by finitely many sets from $\mathcal{G}$. Continuing, we get a decreasing sequence $A_1 \supset A_2 \supset \cdots \supset A_n \supset \ldots$, where $A_n$ is closed, nonempty (infinite, actually), has diam$(A_n \leq 1/n)$, and cannot be covered by finitely many sets from $\mathcal{G}$.

Now here is the problem! Let $x \in \bigcap_{n=1}^\infty A_n$ ($\neq \emptyset$ because $M$ is complete). Then, $x \in G \in \mathcal{G}$ for some $G$ (since $\mathcal{G}$ is an open cover) and so, since $G$ is open, $x \in B_\varepsilon(x) \subset G$ for some $\varepsilon > 0$. But for any $n$ with $1/n < \varepsilon$ we would then have $x \in A_n \subset B_\varepsilon(x) \subset G$. That is, $A_n$ is covered by a single set from $\mathcal{G}$. This is the contradiction that we were looking for.          $\square$

Just look at the tidy form that the nested set theorem takes on in a compact space:

**Corollary 16.** *M is compact if and only if every decreasing sequence of nonempty closed sets has nonempty intersection; that is, if and only if, whenever $F_1 \supset F_2 \supset \ldots$ is a sequence of nonempty closed sets in M, we have $\bigcap_{n=1}^\infty F_n \neq \emptyset$.*

*Proof.*  ($\Rightarrow$) The forward implication is clear from Theorem 50.

($\Leftarrow$) Suppose that every nested sequence of nonempty closed sets in $M$ has nonempty intersection, and let $\{x_n\}$ be a sequence in $M$. Then there is some point $x$ in the nonempty set $\bigcap_{n=1}^\infty \overline{\{x_k \mid k \geq n\}}$ (make sure you see why). It then follows that some subsequence of $\{x_n\}$ must converge to $x$.          $\square$

Note that we no longer need to assume that the diameters of the sets $F_n$ tend to zero; hence, $\bigcap_{n=1}^\infty F_n$ may contain more than one point.

**Corollary 17.** *M is compact if and only if every countable open cover admits a finite subcover. (why?)*

## 1.9 Basics of Measure Theory

We start off this section with an onslaught of definitions. We will need the language presented here as we move forward into the theory of measures.

**Definition 32.** *The **symmetric difference** between two sets E and F (denoted $E \triangle F$) consists of those points that belong to only one of the two sets E or F; that is, $E \triangle F = (E \smallsetminus F) \cup (F \smallsetminus E)$.*

**Definition 33.** *A nonempty collection of sets $\mathfrak{R}$ is called a **ring of sets** if $A \triangle B \in \mathfrak{R}$ and $A \cap B \in \mathfrak{R}$ whenever $A, B \in \mathfrak{R}$. Since*

$$A \cup B = (A \triangle B) \triangle (A \cap B) \qquad \text{and} \qquad A \setminus B = A \triangle (A \cap B),$$

*we also have $A \cup B \in \mathfrak{R}$ and $A \smallsetminus B \in \mathfrak{R}$ whenever $A, B \in \mathfrak{R}$.*

Thus, a ring of sets is a system closed under the operations of taking unions, intersections, differences, and symmetric differences. Clearly, a ring of sets is also closed under the operations of taking finite unions and intersections:

$$\bigcup_{k=1}^{n} A_k \qquad \text{and} \qquad \bigcup_{k=1}^{n} A_k.$$

Note that a ring of sets must also contain the empty set, since $A \smallsetminus A = \emptyset$.

**Definition 34.** *A set E is called the **unit** of a system of sets $\mathcal{S}$ if $E \in \mathcal{S}$ and $A \cap E = A$ for every set $A \in \mathcal{S}$. Clearly E is unique. Thus the unit of $\mathcal{S}$ is just the maximal set of $\mathcal{S}$, i.e. the set containing all other sets of $\mathcal{S}$. A ring of sets with a unit is called an **algebra of sets**.*

A ring of sets is called a *$\sigma$-ring* if it contains the union $\bigcup_{k=1}^{\infty} A_k$ whenever it contains the sets $A_1, A_2, \ldots, A_k, \ldots$. Furthermore, a $\sigma$-ring with a unit is called a *$\sigma$-algebra*. In other words:

**Definition 35.** *Let X be a set. An **algebra** is a collection $\mathcal{A}$ of subsets of X such that:*

i) *$\emptyset \in \mathcal{A}$ and $X \in \mathcal{A}$;*

ii) *if $A \in \mathcal{A}$, then $A^c \in \mathcal{A}$;*

iii) *if $A_1, \ldots, A_n \in \mathcal{A}$, then $\bigcup_{i=1}^{n} A_i$ and $\bigcap_{i=1}^{n} A_i$ are in $\mathcal{A}$.*

*We say that the algebra $\mathcal{A}$ is a **$\sigma$-algebra** if in addition:*

iv) *whenever $A_1, A_2, \cdots \in \mathcal{A}$, then $\bigcup_{i=1}^{\infty} A_i$ and $\bigcap_{i=1}^{\infty} A_i$ are in $\mathcal{A}$ (both the union and intersection must be countable).*

In other words, a $\sigma$-algebra is an algebra of sets, completed to include countably infinite operations; that is, a ***$\sigma$-algebra of sets*** is a collection of subsets that is closed under countable unions, countable intersections, and complements.

**Example 28.**    *a) If X is any set, $\mathcal{P}(X)$ and $\{\emptyset, X\}$ are $\sigma$-algebras.*

b) *If X is uncountable, then*

$$\mathcal{A} = \{E \subset X \mid E \text{ is countable or } E^c \text{ is countable}\}$$

*is a $\sigma$-algebra, called the **$\sigma$-algebra of countable or co-countable sets**. (The point here is that if $\{E_j\} \subset \mathcal{A}$, then $\bigcup_{j=1}^{\infty} E_j$ is countable if all $E_j$ are countable and is co-countable otherwise.)*

**Definition 36.** *Let X be equipped with a $\sigma$-algebra $\mathcal{A}$. A **measure** on $\mathcal{A}$ (or on $(X, \mathcal{A})$, or simply on X if $\mathcal{A}$ is understood) is a function $\mu \colon \mathcal{A} \to [0, \infty]$ such that:*

i) *$\mu(\emptyset) = 0$;*

ii) *If $\{E_j\}$ is a sequence of disjoint sets in $\mathcal{A}$, then $\mu\left(\bigcup_{j=1}^{\infty} E_j\right) = \sum_{j=1}^{\infty} \mu\left(E_j\right)$.*

iii) *Property ii), called **countable additivity**, implies finite additivity: if $E_1, \ldots, E_n$ are disjoint sets in $\mathcal{A}$, then $\mu\left(\bigcup_{j=1}^{n} E_j\right) = \sum_{j=1}^{n} \mu\left(E_j\right)$ (because one can take $E_j = \emptyset$ for $j > n$).*

A function $\mu$ that satisfies i) and iii) but not necessarily ii) is called a ***finitely additive measure***. A measure is obviously also a finitely additive measure, but the converse is not true (a finitely additive measure doesn't satisfy all three properties; hence it's not a measure).

**Definition 37.** *The pair $(X, \mathcal{A})$ is called a **measurable space** and the sets in $\mathcal{A}$ are called **measurable sets**, i.e., a set $A$ is **measurable** (or $\mathcal{A}$-**measurable**) if $A \in \mathcal{A}$. If $\mu$ is a measure on $(X, \mathcal{A})$, then $(X, \mathcal{A}, \mu)$ is called a **measure space**.*

**Definition 38.** *A measure space $X$ is said to be $\sigma$-**finite** if $X$ can be written as the union of countably many measurable sets of finite measure; that is, if $X = \bigcup_i X_i$ and $m(X_i) < \infty$ for all $i$, each $X_i$ being measurable.*

### 1.9.1 Exterior Measure

**Definition 39.** *If $E$ is any subset of $\mathbb{R}^d$, the **exterior measure** (also known as **outer measure**) of $E$, denoted $m_*(E)$, is given by*

$$m_*(E) = \inf \sum_{n=1}^{\infty} |Q_n|,$$

*where the infimum is taken over all countable coverings $E \subset \bigcup_{n=1}^{\infty} Q_n$ by closed cubes.*

The exterior measure is always non-negative but could be infinite, so that in general we have $0 \leq m_*(E) \leq \infty$, and therefore takes values in the extended positive numbers. It follows immediately from the definition of $m_*$ that for every $\varepsilon > 0$, there exists a covering $E \subset \bigcup_{n=1}^{\infty} Q_n$ with

$$\sum_{n=1}^{\infty} m_*(Q_n) \leq m_*(E) + \varepsilon.$$

The relevant properties of exterior measure are now listed in a series of observations:

- Observation 1: (MONOTONICITY) If $E_1 \subset E_2$, then $m_*(E_1) \leq m_*(E_2)$.

- Observation 2: (COUNTABLE SUB-ADDITIVITY) If $E = \bigcup_{n=1}^{\infty} E_n$, then

$$m_*(E) \leq \sum_{n=1}^{\infty} m_*(E_n).$$

- Observation 3: If $E \subset \mathbb{R}^d$, then $m_*(E) = \inf m_*(\mathcal{O})$, where the infimum is taken over all open sets $\mathcal{O}$ containing $E$.

- Observation 4: If $E = E_1 \cup E_2$ and $d(E_1, E_2) > 0$, then $m_*(E) = m_*(E_1) + m_*(E_2)$.

- Observation 5: If a set $E$ is the countable union of almost disjoint cubes $E = \bigcup_{n=1}^{\infty} Q_n$, then $m_*(E) = \sum_{n=1}^{\infty} |Q_n|$.

**Theorem 51** (CARATHÉODORY'S THEOREM). *If $m_*$ is an outer measure on X, the collection $\mathcal{A}$ of $m_*$-measurable sets is a $\sigma$-algebra, and the restriction of $m_*$ to $\mathcal{A}$ is a complete measure.*

*Proof.* See [Folland, 2007, p. 29] for a detailed proof. □

## 1.9.2  Borel Sets & Lebesgue Measure

**Definition 40.** *A subset $E \subset \mathbb{R}$ is said to have **measure zero** if for every $\varepsilon > 0$, there exists a countable family of open intervals $\{I_k\}$ such that*

  *i)* $E \subset \bigcup_{k=1}^{\infty} I_k$;

  *ii)* $\sum_{k=1}^{\infty} |I_k| < \varepsilon$, where $|I_k|$ denotes the length of the interval $I_k$.

The first condition says that the union of the intervals covers $E$, and the second that this union is arbitrarily small. It follows from this that any finite set of points has measure 0. As a matter of fact, it is true that a countable set of points has measure 0, even though the proof of this argument requires a more subtle treatment.

**Definition 41.** *A subset $E$ of $\mathbb{R}^d$ is said to be **Lebesgue measurable** (or simply **measurable**), if for any $\varepsilon > 0$ there exists an open set $\mathcal{O}$ with $E \subset \mathcal{O}$ that satisfies $m_*(\mathcal{O} \smallsetminus E) \leq \varepsilon$. (Alternatively, we can say that a set $E$ is measurable if for all $\varepsilon > 0$, there exists a closed set $F$, such that $F \subset E$ and $m_*(E \smallsetminus F) \leq \varepsilon$.)*

Here's yet another way to define a measurable set (this is actually a more widely used definition):

**Definition 42.** *A set E is **measurable** if for all $A \subset \mathbb{R}^d$, we have*

$$m_*(A) = m_*(E \cap A) + m_*\left(E^c \cap A\right).$$

**Definition 43.** *If E is measurable, we define its **Lebesgue measure** (or simply **measure**) by $m(E) = m_*(E)$. That is, if E is measurable, then its measure is the same as its outer measure.*

Clearly, the Lebesgue measure inherits all the features contained in Observations $1 - 5$ of the exterior measure. Immediately from the definition we find the following six properties:

Property 1: Every open set in $\mathbb{R}^d$ is measurable.

Property 2: If $m_*(E) = 0$, then $E$ is measurable. In particular, if $F$ is a subset of a set of exterior measure 0, then $F$ is measurable.

Property 3: A countable union of measurable sets is measurable.

Property 4: Closed sets are measurable.

Property 5: The complement of a measurable set is measurable.

Property 6: A countable intersection of measurable sets is measurable.

**Theorem 52.** *If $E_1, E_2, \ldots$ are disjoint measurable sets, and $E = \bigcup_{j=1}^{\infty} E_j$, then $m(E) = \sum_{j=1}^{\infty} m\left(E_j\right)$.*

*Proof.* First, we assume further that each $E_j$ is bounded. Then, for each $j$, by applying the definition of measurability to $E_j^c$, we can choose a closed subset $F_j \subset E_j$ with $m_*\left(E_j \smallsetminus F_j\right) \leq \varepsilon/2^j$. For each fixed $N$, the sets $F_1, \ldots, F_N$ are compact and disjoint, so that $m\left(\bigcup_{j=1}^{N} F_j\right) = \sum_{j=1}^{N} m\left(F_j\right)$. Since $\bigcup_{j=1}^{N} F_j \subset E$, we must have

$$m(E) \geq \sum_{j=1}^{N} m\left(F_j\right) \geq \sum_{j=1}^{N} m\left(E_j\right) - \varepsilon.$$

Letting $N$ tend to infinity, since $\varepsilon$ was arbitrary, we find that

$$m(E) \geq \sum_{j=1}^{\infty} m\left(E_j\right).$$

Since the reverse inequality always holds (by sub-additivity in Observation 2 above), this concludes the proof when each $E_j$ is bounded.

Now in the general case, we select any sequence of cubes $\{Q_k\}$ that increases to $\mathbb{R}^d$, in the sense that $Q_k \subset Q_{k+1}$ for all $k \geq 1$ and $\bigcup_{k=1}^{\infty} Q_k = \mathbb{R}^d$. We then let $S_1 = Q_1$ and $S_k = Q_k - Q_{k-1}$ for $k \geq 2$. If we define measurable sets by $E_{j,k} = E_j \cap S_k$, then $E = \cup_{j,k} E_{j,k}$. This union is disjoint and every $E_{j,k}$ is bounded. Moreover $E_j = \bigcup_{k=1}^{\infty} E_{j,k}$ and this union is also disjoint. Putting these facts together, and using what has already been proved, we obtain

$$m(E) = \sum_{j,k} m\left(E_{j,k}\right) = \sum_j \sum_k m\left(E_{j,k}\right) = \sum_j m\left(E_j\right),$$

as claimed. □

With this theorem, the countable additivity of the Lebesgue measure on measurable sets has been established. This result provides the necessary connection between the following:

- our primitive notion of volume given by the exterior measure;

- the more refined idea of measurable sets;

- the countably infinite operations allowed on these sets.

**Theorem 53.** *Suppose $E$ is a measurable subset of $\mathbb{R}^d$. Then, for every $\varepsilon > 0$,*

    *i) there exists an open set $\mathcal{O}$ with $E \subset \mathcal{O}$ and $m(\mathcal{O} \smallsetminus E) \leq \varepsilon$;*

    *ii) there exists a closed set $F$ with $F \subset E$ and $m(E \smallsetminus F) \leq \varepsilon$;*

    *iii) if $m(E)$ is finite, there exists a compact set $K$ with $K \subset E$ and $m(E \smallsetminus K) \leq \varepsilon$;*

    *iv) if $m(E)$ is finite, there exists a finite union $F = \bigcup_{n=1}^{N} Q_n$ of closed cubes such that $m(E \triangle F) \leq \varepsilon$.*

Here are some invariance properties of Lebesgue measure:

- A crucial property of Lebesgue measure in $\mathbb{R}^d$ is its ***translation-invariance***, which can be stated as follows: if $E$ is a measurable set and $h \in \mathbb{R}^d$, then the set $E_h = E + h = \{x + h \mid x \in E\}$ is also measurable, and $m(E + h) = m(E)$.

- Similarly, we have the relative ***dilation-invariance*** of Lebesgue measure: suppose $\delta > 0$, and denote by $\delta E$ the set $\{\delta x \mid x \in E\}$. We can then assert that $\delta E$ is measurable whenever $E$ is, and $m(\delta E) = \delta^d m(E)$ (where $d$ is the dimension of $\mathbb{R}^d \supset E$).

- One can also easily see that the Lebesgue measure is ***reflection-invariant***. That is, whenever $E$ is measurable, so is $-E = \{-x \mid x \in E\}$ and $m(-E) = m(E)$.

The collection of all subsets of $\mathbb{R}^d$ is of course a $\sigma$-algebra. A more interesting and relevant example consists of all measurable sets in $\mathbb{R}^d$, which also forms a $\sigma$-algebra. Another $\sigma$-algebra, which plays a vital role in analysis, is the ***Borel $\sigma$-algebra*** in $\mathbb{R}^d$, denoted by $\mathfrak{B}_{\mathbb{R}^d}$, which by definition is the smallest $\sigma$-algebra that contains all open sets. Elements of this $\sigma$-algebra are called ***Borel sets***. Here the term "smallest" means that if $S$ is any $\sigma$-algebra that contains all open sets in $\mathbb{R}^d$, then necessarily $\mathfrak{B}_{\mathbb{R}^d} \subset S$. Since we observe that any intersection (not necessarily countable) of $\sigma$-algebras is again a $\sigma$-algebra, we may define $\mathfrak{B}_{\mathbb{R}^d}$ as the intersection of all $\sigma$-algebras that contain the open sets. This shows the existence and uniqueness of the Borel $\sigma$-algebra.

More generally:

**Definition 44.** *If $X$ is any topological space (c.f. Chapter 2), the $\sigma$-algebra generated by the family of open sets in $X$ (or equivalently, by the closed sets in $X$) is called the **Borel $\sigma$-algebra**, denoted $\mathcal{B}_X$.*

Let's try to list the Borel sets in order of their complexity:

- We start with the open and closed sets, which are the simplest Borel sets.

- Next in order would come countable intersections of open sets; such sets are called $G_\delta$ sets. Alternatively, one could consider their complements, the countable union of closed sets, called the $F_\sigma$ sets.

- We can then consider a countable union of $G_\delta$ sets (called a $G_{\delta\sigma}$ set); a countable intersection of $F_\sigma$ sets (called an $F_{\sigma\delta}$ set); and so forth. . .

Since open sets and closed sets are measurable, we conclude that the Borel $\sigma$-algebra is contained in the $\sigma$-algebra of measurable sets. Naturally, we may ask if this inclusion is

strict: do there exist Lebesgue measurable sets which are not Borel sets? The answer is "yes." To see why we need to define what a complete measure is:

**Definition 45.** *A **complete measure** is a measure whose domain includes all subsets of null sets.*

In other words, from the point of view of the Borel sets, the Lebesgue sets arise as the completion of the $\sigma$-algebra of Borel sets, that is, by adjoining all subsets of Borel sets of measure zero. This is an immediate consequence of the corollary below.

**Corollary 18.** *A subset $E$ of $\mathbb{R}^d$ is measurable*

- *if and only if $E$ differs from a $G_\delta$ by a set of measure zero;*

- *if and only if $E$ differs from an $F_\sigma$ by a set of measure zero.*

Now this is the theorem that justifies our discussion of complete measures:

**Theorem 54.** *Suppose $(X, \mathcal{A}, \mu)$ is a measure space. Let*

$$\mathfrak{N} = \{N \in \mathcal{A} \mid \mu(N) = 0\} \qquad and \qquad (\bar{\mathcal{A}} = \{E \cup F \mid E \in \mathcal{A} \text{ and } F \subset N \text{ for some } N \in \mathfrak{N}\}.$$

*Then $\bar{\mathcal{A}}$ is a $\sigma$-algebra, and there is a unique extension $\bar{\mu}$ of $\mu$ to a complete measure on $\bar{\mathcal{A}}$.*

*Proof.* See [Folland, 2007, p. 27] for a detailed proof. ☐

Now it is clear what was meant by "completion" before; the measure $\bar{\mu}$ is the completion of $\mu$, and $\bar{\mathcal{A}}$ is the completion of $\mathcal{A}$ with respect to $\mu$.

**Definition 46.** *Let $\{X_\alpha\}_{\alpha \in A}$ be an indexed collection of nonempty sets, $X = \prod_{\alpha \in A} X_\alpha$, and $\pi_\alpha \colon X \to X_\alpha$ the coordinate maps. If $M_\alpha$ is a $\sigma$-algebra on $X_\alpha$ for each $\alpha$, then the **product $\sigma$-algebra** on $X$ is the $\sigma$-algebra generated by*

$$\left\{ \pi_\alpha^{-1}(E_\alpha) \mid E_\alpha \in M_\alpha, \alpha \in A \right\}.$$

*We denote this $\sigma$-algebra by $\bigoplus_{\alpha \in A} M_\alpha$. (If $A = \{1, \ldots, n\}$ we also write $\bigoplus_{k=1}^n M_k$ or $M_1 \oplus \cdots \oplus M_n$.)*

**Proposition 16.** *If $A$ is countable, then $\bigoplus_{\alpha \in A} M_\alpha$ is the $\sigma$-algebra generated by $\{\prod_{\alpha \in A} E_\alpha \mid E_\alpha \in M_\alpha\}$.*

**Proposition 17.** *Suppose that $M_\alpha$ is generated by $\mathcal{E}_\alpha$, $\alpha \in A$. Then $\bigoplus_{\alpha \in A} M_\alpha$ is generated by $\mathcal{F}_1 = \{\pi_\alpha^{-1}(E_\alpha) \mid E_\alpha \in \mathcal{E}_\alpha, \alpha \in A\}$. If $A$ is countable and $X_\alpha \in \mathcal{E}_\alpha$ for all $\alpha$, then $\bigoplus_{\alpha \in A} M_\alpha$ is generated by $\mathcal{F}_2 = \{\prod_{\alpha \in A} E_\alpha : E_\alpha \in \mathcal{E}_\alpha\}$.*

**Proposition 18.** *Let $X_1, \ldots, X_n$ be metric spaces and let $X = \prod_{k=1}^{n} X_k$, equipped with the product metric. Then $\bigoplus_{k=1}^{n} \mathcal{B}_{X_k} \subset \mathcal{B}_X$. If the $X_k$'s are separable, then $\bigoplus_{k=1}^{n} \mathcal{B}_{X_k} = \mathcal{B}_X$.*

**Corollary 19.** $\mathcal{B}_{\mathbb{R}^d} = \bigoplus_{k=1}^{n} \mathcal{B}_{\mathbb{R}}$.

### 1.9.3   Construction of a Non-Measurable Set

The construction of a non-measurable (Vitali) set $N$ uses the axiom of choice, and rests on a simple equivalence relation among real numbers in $[0, 1]$:

$$\text{We write } x \sim y \text{ whenever } x - y \in \mathbb{Q}.$$

Note that this is an equivalence relation since it satisfies the reflexive, symmetric, and transitive properties. Since equivalence classes partition a set into distinct cells, we know that two equivalence classes either are disjoint or coincide; thus the interval $[0, 1]$ is the disjoint union of all equivalence classes that live in this interval, that is

$$[0, 1] = \bigcup_\alpha \mathcal{E}_\alpha,$$

where each $\mathcal{E}_\alpha$ represents a unique equivalence class.

Now we construct the (Vitali) set $N$ by choosing exactly one element $x_\alpha$ from each $\mathcal{E}_\alpha$ (this is justified by using the axiom of choice), and setting $N = \{x_\alpha\}$. Here's the important result:

**Theorem 55.** *The Vitali set $N$ constructed above is not measurable.*

*Proof.* Assume, to the contrary, that $N$ is measurable. Let $\{r_k\}$ be an enumeration of all the rationals in $[-1, 1]$, and consider the translates

$$N_k = N + r_k.$$

Note that the sets $N_k$ are disjoint. To see why this is true, suppose that the intersection $N_k \cap N_{k'}$ is nonempty. Then there exist rationals $r_k \neq r_{k'}$ and $\alpha$ and $\beta$ with

$$x_\alpha + r_k = x_\beta + r_{k'}$$

which implies that

$$x_\alpha - x_\beta = r_{k'} - r_k.$$

But this means that $\alpha \neq \beta$ and $x_\alpha - x_\beta$ is rational, which in turn implies that $x_\alpha \sim x_\beta$. This contradicts the fact that $N$ contains only one representative of each equivalence class.

We also claim that

$$[0, 1] \subset \bigcup_{k=1}^{\infty} N_k \subset [-1, 2]. \tag{1.5}$$

To see why, notice that if $x \in [0, 1]$, then $x \sim x_\alpha$ for some $\alpha$, and therefore $x - x_\alpha = r_k$ for some $k$. Hence $x \in N_k$ for some $k$ and the first inclusion holds. The second inclusion above is straightforward since each $N_k$ is contained in $[-1, 2]$ by construction.

Now we may conclude the proof of the theorem. If $N$ were measurable, then so would be $N_k$ for all $k$, and since the union $\bigcup_{k=1}^{\infty} N_k$ is disjoint, the inclusions in (1.5) yield

$$1 \leq \sum_{k=1}^{\infty} m\left(N_k\right) \leq 3.$$

Since $N_k$ is a translate of $N$, we must have $m\left(N_k\right) = m(N)$ for all $k$. Consequently,

$$1 \leq \sum_{k=1}^{\infty} m(N) \leq 3.$$

This is the desired contradiction, since neither $m(N) = 0$ nor $m(N) > 0$ is possible. In other words, $m(N) = 0$ is not possible by the above inequality, and $m(N) > 0$ is not possible either because we are trying to find the measure of a countable set, which would have measure zero if any. $\qquad\qquad\square$

### 1.9.4   Measurable Functions

Our starting point is the notion of a ***characteristic function*** of a set $E$, which is defined by

$$\chi_E(x) = \begin{cases} 1 & \text{if } x \in E, \\ 0 & \text{if } x \notin E. \end{cases}$$

The next step is to pass to the functions that are the building blocks of integration theory. For the Riemann integral it is in effect the class of step functions, with each given as a finite sum

$$f = \sum_{k=1}^{N} a_k \chi_{R_k},$$

where each $R_k$ is a rectangle, and the $a_k$ are constants. However, for the Lebesgue integral we need a more general notion, as we shall see later on.

**Definition 47.** *A **simple function** is a finite sum*

$$f = \sum_{k=1}^{N} a_k \chi_{E_k},$$

*where each $E_k$ is a measurable set of finite measure, and the $a_k$ are constants.*

**Definition 48.** *A function $f$ defined on a measurable subset $E$ of $\mathbb{R}^d$ is said to be a **measurable function**, if for all $a \in \mathbb{R}$, the set*

$$f^{-1}\left([-\infty, a)\right) = \{x \in E \mid f(x) < a\}$$

*is measurable. (To simplify our notation, we shall often denote the set $\{x \in E \mid f(x) < a\}$ simply by $\{f < a\}$ whenever no confusion is possible.)*

    In the same way, one can show that if $f$ is finite-valued, then it is measurable if and only if the sets $\{a < f < b\}$ are measurable for every $a, b \in \mathbb{R}$. Similar conclusions hold for whichever combination of strict or weak inequalities one chooses. For example, if $f$ is finite-valued, then it is measurable if and only if $\{a \leq f \leq b\}$ is measurable for all $a, b \in \mathbb{R}$. By the same arguments one sees the following:

- Property 1: The finite-valued function $f$ is measurable if and only if $f^{-1}(\mathcal{O})$ is measurable for every open set $\mathcal{O}$, and if and only if $f^{-1}(F)$ is measurable for every closed set $F$. (Remark: Note that this property also applies to extended-valued functions, if we make the additional hypothesis that both $f^{-1}(-\infty)$ and $f^{-1}(\infty)$ are measurable sets.)

- Property 2: If $f$ is continuous on $\mathbb{R}^d$, then $f$ is measurable. If $f$ is measurable and finite-valued, and $\Phi$ is continuous, then $\Phi \circ f$ is measurable. (Remark: In fact, $\Phi$ is continuous, so $\Phi^{-1}(-\infty, a)$ is an open set $\mathcal{O}$, and hence $(\Phi \circ f)^{-1}((-\infty, a)) = f^{-1}(\mathcal{O})$ is measurable. It should be noted, however, that in general it is not true that $f \circ \Phi$ is measurable whenever $f$ is measurable and $\Phi$ is continuous.)

- Property 3: Suppose $\{f_n\}$ is a sequence of measurable functions. Then

$$\sup_n f_n(x), \quad \inf_n f_n(x), \quad \limsup_{n\to\infty} f_n(x), \quad \text{and} \quad \liminf_{n\to\infty} f_n(x)$$

are measurable. (Remark: Proving that $\sup_n f_n$ is measurable requires noting that $\{\sup_n f_n > a\} = \bigcup_n \{f_n > a\}$. This also yields the result for $\inf_n f_n$, since this quantity equals $-\sup_n(-f_n(x))$. The result for the $\limsup$ and $\liminf$ also follows from the two observations

$$\limsup_{n\to\infty} f_n = \inf_k \left\{ \sup_{n \geq k} f_n \right\} \quad \text{and} \quad \liminf_{n\to\infty} f_n = \sup_k \left\{ \inf_{n \geq k} f_n \right\}.$$

- Property 4: If $\{f_n\}$ is a collection of measurable functions, and $\lim_{n\to\infty} f_n(x) = f(x)$, then $f$ is measurable. (Remark: Since $f(x) = \limsup_{n\to\infty} f_n(x) = \liminf_{n\to\infty} f_n(x)$, this property is a consequence of property 3.)

- Property 5: If $f$ and $g$ are measurable, then

    i) The integer powers $f^k$, for $k \geq 1$ are measurable.

    ii) $f + g$ and fg are measurable if both $f$ and $g$ are finite-valued.

Remark: For i) we simply note that if $k$ is odd, then $\{f^k > a\} = \{f > a^{1/k}\}$, and if $k$ is even and $a \geq 0$, then

$$\left\{f^k > a\right\} = \left\{f > a^{1/k}\right\} \bigcup \left\{f > -a^{1/k}\right\}.$$

For ii), we first see that $f + g$ is measurable because

$$\{f + g > a\} = \bigcup_{r \in \mathbb{Q}} \{f > a - r\} \bigcap \{g > r\}.$$

Finally, $fg$ is measurable because of the previous results and the fact that

$$fg = \frac{1}{4} \left[ (f + g)^2 - (f - g)^2 \right].$$

**Definition 49.** *We shall say that two functions $f$ and $g$ defined on a set $E$ are **equal almost everywhere**, and write*

$$f(x) = g(x) \quad a.e. \ x \in E$$

*if the set $\{x \in E \mid f(x) \neq g(x)\}$ has measure zero (we sometimes abbreviate this by saying that $f = g$ a.e.). More generally, a property or statement is said to hold almost everywhere (a.e.) if it is true except on a set of measure zero.*

One sees easily that if $f$ is measurable and $f = g$ a.e., then $g$ is measurable. This follows at once from the fact that $\{f < a\}$ and $\{g < a\}$ differ by a set of measure zero. Moreover, all the properties stated above this definition can be relaxed to conditions holding almost everywhere. For instance, if $\{f_n\}$ is a collection of measurable functions, and $\lim_{n \to \infty} f_n(x) = f(x)$ a.e., then $f$ is measurable. This observation gives us a sixth property of measurable functions:

- Property 6: Suppose $f$ is measurable, and $f(x) = g(x)$ for a.e. $x$. Then $g$ is also measurable.

### 1.9.5   Approximation By Simple or Step Functions

The theorems in this section are all of the same nature and provide further insight in the structure of measurable functions. We begin by approximating pointwise, non-negative measurable functions by simple functions:

**Theorem 56.** *Suppose $f$ is a non-negative measurable function on $\mathbb{R}^d$. Then there exists an increasing sequence of non-negative simple functions $\{\varphi_k\}$ that converges pointwise to $f$, namely,*

$$\varphi_k(x) \leq \varphi_{k+1}(x) \qquad and \qquad \lim_{k \to \infty} \varphi_k(x) = f(x) \quad for \ all \ x.$$

*Proof.* We begin first with a truncation. For $N \geq 1$, let $Q_N$ denote the cube centered at the origin and of side length $N$. Then we define

$$F_N(x) = \begin{cases} f(x) & \text{if } x \in Q_N \text{ and } f(x) \leq N, \\ N & \text{if } x \in Q_N \text{ and } f(x) > N, \\ 0 & \text{otherwise.} \end{cases}$$

Then, $F_N(x) \to f(x)$ as $N$ tends to infinity for all $x$. Now, we partition the range of $F_N$, namely $[0, N]$, as follows. For fixed $N, M \geq 1$, we define

$$E_{\ell,M} = \left\{ x \in Q_N : \frac{\ell}{M} < F_N(x) \leq \frac{\ell+1}{M} \right\} \quad \text{for } 0 \leq \ell < NM.$$

Then we may form

$$F_{N,M}(x) = \sum_{\ell} \frac{\ell}{M} \chi_{E_{\ell,M}}(x).$$

Each $F_{N,M}$ is a simple function that satisfies $0 \leq F_N(x) - F_{N,M}(x) \leq 1/M$ for all $x$. If we now choose $N = M = 2^k$ with $k \leq 1$ integral, and let $\varphi_k = F_{2^k,2^k}$, then we see that $0 \leq F_M(x) - \varphi_k(x) \leq 1/2^k$ for all $x$, $\{\varphi_k\}$ is increasing, and this sequence satisfies all the desired properties. $\qquad\square$

Note that the result holds for non-negative functions that are extended-valued, if the limit $+\infty$ is allowed. We now drop the assumption that $f$ is nonnegative, and also allow the extended limit $-\infty$:

**Theorem 57.** *Suppose $f$ is measurable on $\mathbb{R}^d$. Then there exists a sequence of simple functions $\{\varphi_k\}$ that satisfies*

$$|\varphi_k(x)| \leq |\varphi_{k+1}(x)| \qquad \text{and} \qquad \lim_{k \to \infty} \varphi_k(x) = f(x) \qquad \text{for all } x.$$

*In particular, we have $|\varphi_k(x)| \leq |f(x)|$ for all $x$ and $k$.*

*Proof.* We decompose the function $f$ as $f(x) = f^+(x) - f^-(x)$, where

$$f^+(x) = \max(f(x), 0) \quad \text{and} \quad f^{-1}(x) = \max(-f(x), 0).$$

Since both $f^+$ and $f^{-1}$ are non-negative, Theorem 56 yields two increasing sequences of nonnegative simple functions $\left\{\varphi_k^{(1)}(x)\right\}$ and $\left\{\varphi_k^{(2)}(x)\right\}$ which converge pointwise to $f^+$ and $f^{-1}$, respectively. Then, if we let

$$\varphi_k(x) = \varphi_k^{(1)}(x) - \varphi_k^{(2)}(x),$$

we see that $\varphi_k(x)$ converges to $f(x)$ for all $x$. Finally, the sequence $\{|\varphi_k|\}$ is increasing because the definition of $f^+$, $f^-$ and the properties of $\varphi_k^{(1)}$ and $\varphi_k^{(2)}$ imply that

$$|\varphi_k(x)| = \varphi_k^{(1)}(x) + \varphi_k^{(2)}(x). \qquad \square$$

We may now go one step further, and approximate by step functions. Here, in general, the convergence may hold only almost everywhere:

**Theorem 58.** *Suppose $f$ is measurable on $\mathbb{R}^d$. Then there exists a sequence of step functions $\{\psi_k\}$ that converges pointwise to $f(x)$ for almost every $x$.*

*Proof.* By the previous result, it suffices to show that if $E$ is a measurable set with finite measure, then $f = \chi_E$ can be approximated by step functions. To this end, we recall that for every $\varepsilon$ there exist cubes $Q_1, \ldots, Q_N$ such that

$$m\left(E \triangle \bigcup_{j=1}^{N} Q_j\right) \leq \varepsilon.$$

By considering the grid formed by extending the sides of these cubes, we see that there exist almost disjoint rectangles $\widetilde{R}_1, \ldots, \widetilde{R}_M$ such that

$$\bigcup_{j=1}^{N} Q_j = \bigcup_{j=1}^{M} \widetilde{R}_j.$$

By taking rectangles $R_j$ contained in $\widetilde{R}_j$, and slightly smaller in size, we find a collection of disjoint rectangles that satisfy

$$m\left(E \triangle \bigcup_{j=1}^{M} R_j\right) \leq 2\varepsilon.$$

Therefore

$$f(x) = \sum_{j=1}^{M} \chi_{R_j}(x),$$

except possibly on a set of measure $\leq 2\varepsilon$. Consequently, for every $k \geq 1$, there exists a step function $\psi_k(x)$ such that if

$$E_k = \{x \mid f(x) \neq \psi_k(x)\},$$

then $m(E_k) \leq 2^{-k}$. If we let $F_k = \bigcup_{j=K+1}^{\infty} E_j$ and $F = \bigcap_{K=1}^{\infty} F_K$, then $m(F) = 0$ since $m(F_K) \leq 2^{-K}$, and $\psi_k(x) \to f(x)$ for all $x$ in the complement of $F$, which is the desired result.     □


Although the notions of measurable sets and measurable functions represent new tools, we should not overlook their relation to the older concepts they replaced. The British mathematician John Edensor Littlewood aptly summarized these connections in the form of three principles that provide a useful intuitive guide in the initial study of the theory. Here are the famous *Littlewood's three principles*:

L1)  Every set is nearly a finite union of intervals.

L2)  Every function is nearly continuous.

L3)  Every convergent sequence is nearly uniformly convergent.

The sets and functions referred to above are of course assumed to be measurable. The catch is in the word "nearly," which has to be understood appropriately in each context.

**Theorem 59** (EGOROV'S THEOREM). *Suppose $\{f_k\}$ is a sequence of measurable functions defined on a measurable set $E$ with $m(E) < \infty$, and assume that $f_k \to f$ a.e on $E$. Given $\varepsilon > 0$, we can find a closed set $A_\varepsilon \subset E$ such that $m(E \setminus A_\varepsilon) \leq \varepsilon$ and $f_k \to f$ uniformly on $A_\varepsilon$.*

*Proof.* We may assume WLOG that $f_k(x) \to f(x)$ for every $x \in E$. For each pair of nonnegative integers $n$ and $k$, let

$$E_k^n = \left\{x \in E : \left|f_j(x) - f(x)\right| < 1/n, \text{ for all } j > k\right\}.$$

Now fix $n$ and note that $E_k^n \subset E_{k+1}^n$, and $E_k^n \nearrow E$ as $k$ tends to infinity. By a previous corollary, we find that there exists $k_n$ such that $m(E \setminus E_{k_n}^n) < 1/2^n$. By construction, we then have

$$\left|f_j(x) - f(x)\right| < 1/n \qquad \text{whenever } j > k_n \text{ and } x \in E_{k_n}^n.$$

We choose $N$ so that $\sum_{n=N}^{\infty} 2^{-n} < \varepsilon/2$, and let

$$\widetilde{A}_{\varepsilon} = \bigcap_{n \geq N} E_{k_n}^{n}.$$

We first observe that

$$m\left(E \smallsetminus \widetilde{A}_{\varepsilon}\right) \leq \sum_{n=N}^{\infty} m\left(E \smallsetminus E_{k_n}^{n}\right) < \frac{\varepsilon}{2}.$$

Next, if $\delta > 0$, we choose $n \geq N$ such that $1/n < \delta$, and note that $x \in \widetilde{A}_{\varepsilon}$ implies $x \in E_{k_n}^{n}$. We see therefore that $\left|f_j(x) - f(x)\right| < \delta$ whenever $j > k_n$. Hence $f_k$ converges uniformly to $f$ on $\widetilde{A}_{\varepsilon}$. Finally, by a previous theorem we can choose a closed subset $A_{\varepsilon} \subset \widetilde{A}_{\varepsilon}$ with $m\left(\widetilde{A}_{\varepsilon} \smallsetminus A_{\varepsilon}\right) < \varepsilon/2$. As a result, we have $m\left(E \smallsetminus A_{\varepsilon}\right) < \varepsilon$ and the theorem is proved. $\qquad\square$

The next theorem attests to the validity of the second of Littlewood's principles:

**Theorem 60** (LUSIN'S THEOREM). *Suppose $f$ is measurable and finite valued on a set $E$ of finite measure. Then for every $\varepsilon > 0$ there exists a closed set $F_{\varepsilon}$, with*

$$F_{\varepsilon} \subset E \qquad and \qquad m\left(E \smallsetminus F_{\varepsilon}\right) \leq \varepsilon$$

*and such that $f|_{F_{\varepsilon}}$ is continuous.*

*Proof.* Let $f_n$ be a sequence of step functions so that $f_n \to f$ a.e. Then we may find sets $E_n$ so that $m\left(E_n\right) < 1/2^n$ and $f_n$ is continuous outside $E_n$. By Egorov's theorem, we may find a set $A_{\varepsilon/3}$ on which $f_n \to f$ uniformly and $m\left(E \smallsetminus A_{\varepsilon/3}\right) \leq \varepsilon/3$. Then we consider

$$F' = A_{\varepsilon/3} \setminus \bigcup_{n \geq N} E_n$$

for $N$ so large that $\sum_{n \geq N} 1/2^n < \varepsilon/3$. Now for every $n \geq N$ the function $f_n$ is continuous on $F'$; thus $f$ (being the uniform limit of $\{f_n\}$ by Egorov's theorem) is also continuous on $F'$. Now to finish the proof we merely need to approximate the set $F'$ by a closed set $F_{\varepsilon} \subset F'$ such that $m\left(F' \smallsetminus F_{\varepsilon}\right) < \varepsilon/3$. $\qquad\square$

**Remark:** The conclusion of the theorem states that if $f$ is viewed as a function defined only on $F_{\varepsilon}$, then $f$ is continuous. However, the theorem does not make the stronger assertion that the function $f$ defined on $E$ is continuous at the points of $F_{\varepsilon}$.

## 1.10 Lebesgue Integration

The general notion of the Lebesgue integral on $\mathbb{R}^d$ will be defined in a step-by-step fashion, proceeding successively to increasingly larger families of functions. At each stage we shall see that the integral satisfies elementary properties such as linearity and monotonicity, and we show appropriate convergence theorems that amount to interchanging the integral with limits. At the end of the process we shall have achieved a general theory of integration that will be decisive in the study of further problems. We emphasize from the onset that all functions are assumed to be measurable.

### 1.10.1   Stage One: Simple Functions

Recall that a simple function $\varphi$ is a finite sum

$$\varphi(x) = \sum_{k=1}^{N} \alpha_k \chi_{E_k}(x), \tag{1.6}$$

where the $E_k$ are measurable sets of finite measure and the $\alpha_k$ are constants. A complication that arises from this definition is that a simple function can be written in a multitude of ways as such finite linear combinations; for example, $0 = \chi_E - \chi_E$ for any measurable set $E$ of finite measure. Fortunately, there is an unambiguous choice for the representation of a simple function, which is natural and useful in applications: The ***canonical form of $\varphi$*** is the unique decomposition as in (1.6), where the numbers $\alpha_k$ are distinct and nonzero, and the sets $E_k$ are disjoint.

Finding the canonical form of $\varphi$ is straightforward: Since $\varphi$ can take only finitely many distinct and non-zero values, say $c_1, \ldots, c_M$, we may set

$$F_k = \{ x \mid \varphi(x) = c_k \},$$

and note that the sets $F_k$ are disjoint. Therefore

$$\varphi = \sum_{k=1}^{M} c_k \chi_{F_k}$$

is the desired canonical form of $\varphi$.

**Definition 50.** *If $\varphi$ is a simple function with canonical form $\varphi = \sum_{k=1}^{M} c_k \chi_{F_k}$, then we define the* ***Lebesgue integral*** *of $\varphi$ by*

$$\int_{\mathbb{R}^d} \varphi(x)\,\mathrm{d}x = \int_{\mathbb{R}^d} \sum_{k=1}^{M} c_k \chi_{F_k}(x)\,\mathrm{d}x = \sum_{k=1}^{M} c_k m(F_k).$$

*If $E$ is a measurable subset of $\mathbb{R}^d$ with finite measure, then $\varphi(x)\chi_E(x)$ is also a simple function, and we define*

$$\int_E \varphi(x)\mathrm{d}x = \int \varphi(x)\chi_E(x)\mathrm{d}x.$$

**Remark:** To emphasize the choice of the Lebesgue measure $m$ in the definition of the integral, one sometimes writes $\int_{\mathbb{R}^d} \varphi(x)\mathrm{d}m(x)$ for the Lebesgue integral of $\varphi$. However, as a matter of convenience, we shall often write $\int \varphi(x)\,\mathrm{d}x$ or simply $\int \varphi$ for the integral of $\varphi$ over $\mathbb{R}^d$.

**Proposition 19.** *The integral of simple functions defined above satisfies the following properties:*

*i)* (INDEPENDENCE OF REPRESENTATION) *If $\varphi = \sum_{k=1}^{N} \alpha_k \chi_{E_k}$ is any representation of $\varphi$, then*

$$\int \varphi = \sum_{k=1}^{N} \alpha_k m(E_k).$$

*ii)* (LINEARITY) *If $\varphi$ and $\psi$ are simple, and $\alpha, \beta \in \mathbb{R}$, then*

$$\int (\alpha\varphi + \beta\psi) = \alpha \int \varphi + \beta \int \psi.$$

*iii)* (ADDITIVITY) *If $E$ and $F$ are disjoint subsets of $\mathbb{R}^d$ with finite measure, then*

$$\int_{E \cup F} \varphi = \int_E \varphi + \int_F \varphi.$$

*iv)* (MONOTONICITY) *If $\varphi \leq \psi$ are simple, then*

$$\int \varphi \leq \int \psi.$$

*v)* (TRIANGLE INEQUALITY) *If $\varphi$ is a simple function, then so is $|\varphi|$, and*

$$\left| \int \varphi \right| \leq \int |\varphi|.$$

A proof of this proposition can be found on [Stein, 2005, Pg. 51].

## 1.10.2   Stage Two: Bounded Functions Supported on a Set of Finite Measure

Recall that the *support* of a measurable function $f$ is defined to be the set of all points where $f$ does not vanish; that is,

$$\text{supp} f = \{x \mid f(x) \neq 0\}.$$

We shall also say that $f$ is *supported* on a set $E$, if $f(x) = 0$ whenever $x \notin E$. It follows that, since $f$ is measurable, so is the set supp $f$. We shall next be interested in those bounded measurable functions that have $m(\text{supp} f) < \infty$.

An important result is that if $f$ is a function bounded by $M$ and supported on a set $E$, then there exists a sequence $\{\varphi_n\}$ of simple functions, with each $\varphi_n$ bounded by $M$ and supported on $E$, and such that

$$\varphi_n(x) \to f(x) \qquad \forall x.$$

The key lemma that follows allows us to define the integral for the class of bounded functions supported on sets of finite measure:

**Lemma 11.** *Let $f$ be a bounded function supported on a set $E$ of finite measure. If $\{\varphi_n\}_{n=1}^{\infty}$ is any sequence of simple functions bounded by $M$, supported on $E$, and with $\varphi_n(x) \to f(x)$ for a.e. $x$, then:*

*i)* $\lim_{n \to \infty} \int \varphi_n$ *exists.*

*ii)* *if $f = 0$ a.e., then $\lim_{n \to \infty} \int \varphi_n = 0$.*

*Proof.* The assertions of the lemma would be nearly obvious if we had that $\varphi_n$ converges to $f$ uniformly on $E$. Instead, we recall one of Littlewood's principles, which states that the convergence of a sequence of measurable functions is "nearly" uniform. The precise statement lying behind this principle is Egorov's theorem, which we apply here.

For part i), since the measure of $E$ is finite, given $\varepsilon > 0$ Egorov's theorem guarantees the existence of a (closed) measurable subset $A_\varepsilon \subset E$ such that $m\left(E \setminus A_\varepsilon\right) \leq \varepsilon$, and $\varphi_n \to f$ uniformly on $A_\varepsilon$. Therefore, setting $I_n = \int \varphi_n$ we have that

$$
\begin{aligned}
|I_n - I_m| &\leq \int_E |\varphi_n(x) - \varphi_m(x)|\, dx \\
&= \int_{A_\varepsilon} |\varphi_n(x) - \varphi_m(x)|\, dx + \int_{E \setminus A_\varepsilon} |\varphi_n(x) - \varphi_m(x)|\, dx \\
&\leq \int_{A_\varepsilon} |\varphi_n(x) - \varphi_m(x)|\, dx + 2Mm\left(E \setminus A_\varepsilon\right) \\
&\leq \int_{A_\varepsilon} |\varphi_n(x) - \varphi_m(x)|\, dx + 2M\varepsilon.
\end{aligned}
$$

By the uniform convergence, one has, for all $x \in A_\varepsilon$ and all large $n$ and $m$, the estimate $|\varphi_n(x) - \varphi_m(x)| < \varepsilon$, so we deduce that

$$
|I_n - I_m| \leq m(E)\varepsilon + 2M\varepsilon \qquad \text{for all large } n \text{ and } m.
$$

Since $\varepsilon$ is arbitrary and $m(E) < \infty$, this proves that $\{I_n\}$ is a Cauchy sequence and hence converges, as desired.

To show part ii), we note that if $f = 0$, we may repeat the argument above to find that $|I_n| \leq m(E)\varepsilon + M\varepsilon$, which yields $\lim_{n\to\infty} I_n = 0$, as was to be shown. $\qquad\square$

Using the above lemma we can now turn to the integration of bounded functions that are supported on sets of finite measure. For such a function $f$ we define its *Lebesgue integral* by

$$
\int f(x)\, dx = \lim_{n\to\infty} \int \varphi_n(x)\, dx,
$$

where $\{\varphi_n\}$ is any sequence of simple functions satisfying:

    i) $|\varphi_n| \leq M$.

ii) each $\varphi_n$ is supported on the support of $f$.

iii) $\varphi_n(x) \to f(x)$ for a.e. $x$ as $n$ tends to infinity (we know by Lemma 11 above that this limit exists).

Next, we must first show that $\int f$ is independent of the limiting sequence $\{\varphi_n\}$ used, in order for the integral to be well defined. Therefore, suppose that $\{\psi_n\}$ is another sequence of simple functions that is bounded by $M$, supported on $\operatorname{supp} f$, and such that $\psi_n(x) \to f(x)$ for a.e. $x$ as $n$ tends to infinity. Then, if $\eta_n = \varphi_n - \psi_n$, the sequence $\{\eta_n\}$ consists of simple functions bounded by $2M$, supported on a set of finite measure, and such that $\eta_n \to 0$ a.e. as $n$ tends to infinity. We may therefore conclude, by the second part of Lemma 11, that $\int \eta_n \to 0$ as $n$ tends to infinity. Consequently, the two limits

$$\lim_{n\to\infty} \int \varphi_n(x)\, dx \qquad \text{and} \qquad \lim_{n\to\infty} \int \psi_n(x)\, dx$$

(which exist by the lemma) are indeed equal.

Now, if $E$ is a subset of $\mathbb{R}^d$ with finite measure, and $f$ is bounded with $m(\operatorname{supp} f) < \infty$, then it is natural to define

$$\int_E f(x)\, dx = \int f(x)\chi_E(x)\, dx.$$

Clearly, if $f$ is itself simple, then $\int f$ as defined above coincides with the integral of simple functions presented earlier. This extension of the definition of integration also satisfies all the basic properties of the integral of simple functions:

**Proposition 20.** *Suppose $f$ and $g$ are bounded functions supported on sets of finite measure. Then the following properties hold:*

  *i)* (LINEARITY) *If $a, b \in \mathbb{R}$, then*

$$\int (\alpha f + \beta g) = \alpha \int f + \beta \int g.$$

  *ii)* (ADDITIVITY) *If $E$ and $F$ are disjoint subsets of $\mathbb{R}^d$, then*

$$\int_{E \cup F} f = \int_E f + \int_F f.$$

*iii)* (MONOTONICITY) *If $f \leq g$, then*

$$\int f \leq \int g.$$

*iv)* (TRIANGLE INEQUALITY) *$|f|$ is also bounded, supported on a set of finite measure, and*

$$\left| \int f \right| \leq \int |f|.$$

We are now in a position to prove the first important convergence theorem:

**Theorem 61 (Bounded Convergence Theorem).** *Suppose that $\{f_n\}$ is a sequence of measurable functions that are all bounded by $M$, are supported on a set $E$ of finite measure, and $f_n(x) \to f(x)$ a.e. $x$ as $n \to \infty$. Then $f$ is measurable, bounded, supported on $E$ for a.e. $x$, and*

$$\int |f_n - f| \to 0 \quad \text{as } n \to \infty.$$

*Consequently,*

$$\int f_n \to \int f \quad \text{as } n \to \infty.$$

*Proof.* From the assumptions one sees at once that $f$ is bounded by $M$ almost everywhere and vanishes outside $E$, except possibly on a set of measure zero. Clearly, the triangle inequality for the integral implies that it suffices to prove that $\int |f_n - f| \to 0$ as $n$ approaches infinity. Given $\varepsilon > 0$, we may find, by Egorov's theorem, a measurable subset $A_\varepsilon \subset E$ such that $m(E \smallsetminus A_\varepsilon) \leq \varepsilon$ and $f_n \to f$ uniformly on $A_\varepsilon$. Then, we know that for all sufficiently large $n$ and for all $x \in A_\varepsilon$, we have

$$|f_n(x) - f(x)| \leq \varepsilon.$$

Putting these facts together yields

$$\int_E |f_n(x) - f(x)| \, dx = \int_{A_\varepsilon} |f_n(x) - f(x)| \, dx + \int_{E \smallsetminus A_\varepsilon} |f_n(x) - f(x)| \, dx \tag{1.7}$$

$$\leq \varepsilon \, m(E) + 2M \, m(E \smallsetminus A_\varepsilon) \tag{1.8}$$

for all large $n$. Since $\varepsilon$ is arbitrary, the proof of the theorem is complete. $\qquad \square$

**Remark:** We note that the above convergence theorem is a statement about the interchange of an integral and a limit, since its conclusion simply says that $\lim_{n\to\infty} \int f_n = \int \lim_{n\to\infty} f_n$.

A useful observation that we can make at this point is the following: if $f \geq 0$ is bounded and supported on a set of finite measure $E$ and $\int f = 0$, then $f = 0$ almost everywhere. Indeed, if for each integer $k \geq 1$ we set

$$E_k = \{x \in E \mid f(x) \geq 1/k\},$$

then the fact that $k^{-1}\chi_{E_k}(x) \leq f(x)$ implies

$$k^{-1}m(E_k) \leq \int f,$$

by monotonicity of the integral. Thus $m(E_k) = 0$ for all $k$, and since

$$\{x \mid f(x) > 0\} = \bigcup_{k=1}^{\infty} E_k,$$

we see that $f = 0$ almost everywhere.

To close this section, we briefly return to Riemann integrable functions with the following important result:

**Theorem 62.** *Suppose $f$ is Riemann integrable on the closed interval $[a, b]$. Then $f$ is measurable, and*

$$\int_{[a,b]}^{\mathcal{R}} f(x)\, dx = \int_{[a,b]}^{\mathcal{L}} f(x)\, dx,$$

*where the integral on the left-hand side is the standard Riemann integral, and that on the right-hand side is the Lebesgue integral.*

*Proof.* By definition, a Riemann integrable function is bounded, say $|f(x)| \leq M$, so we need to prove that $f$ is measurable, and then establish the equality of integrals. Again, by definition of Riemann integrability, we may construct two sequences of step functions $\{\varphi_k\}$ and $\{\psi_k\}$ that satisfy the following properties:

(∗) $|\varphi_k(x)| \leq M$ and $|\psi_k(x)| \leq M$ for all $x \in [a, b]$ and $k \geq 1$.

$(**)$ $\varphi_1(x) \le \varphi_2(x) \le \cdots \le f \le \psi_2(x) \le \psi_1(x)$.

$(***)$ $\lim_{k \to \infty} \int_{[a,b]}^{\mathcal{R}} \varphi_k(x)\,\mathrm{d}x = \lim_{k \to \infty} \int_{[a,b]}^{\mathcal{R}} \psi_k(x)\,\mathrm{d}x = \int_{[a,b]}^{\mathcal{R}} f(x)\,\mathrm{d}x$.

Several observations are in order. First, it follows immediately from their definition that for step functions the Riemann and Lebesgue integrals agree; therefore

$$\int_{[a,b]}^{\mathcal{R}} \varphi_k(x)\,\mathrm{d}x = \int_{[a,b]}^{\mathcal{L}} \varphi_k(x)\,\mathrm{d}x \qquad \text{and} \qquad \int_{[a,b]}^{\mathcal{R}} \psi_k(x)\,\mathrm{d}x = \int_{[a,b]}^{\mathcal{L}} \psi_k(x)\,\mathrm{d}x, \qquad (1.9)$$

for all $k \ge 1$.

Next, if we let

$$\widetilde{\varphi}(x) = \lim_{k \to \infty} \varphi_k(x) \qquad \text{and} \qquad \widetilde{\psi}(x) = \lim_{k \to \infty} \psi_k(x),$$

we have $\widetilde{\varphi} < f < \widetilde{\psi}$. Moreover, both $\widetilde{\varphi}$ and $\widetilde{\psi}$ are measurable (being the limit of step functions), and the bounded convergence theorem yields

$$\lim_{k \to \infty} \int_{[a,b]}^{\mathcal{L}} \varphi_k(x)\,\mathrm{d}x = \int_{[a,b]}^{\mathcal{L}} \widetilde{\varphi}(x)\,\mathrm{d}x$$

$$\lim_{k \to \infty} \int_{[a,b]}^{\mathcal{L}} \psi_k(x)\,\mathrm{d}x = \int_{[a,b]}^{\mathcal{L}} \widetilde{\psi}(x)\,\mathrm{d}x.$$

This result, together with property $(***)$ and equations (1.9), yield

$$\int_{[a,b]}^{\mathcal{L}} \left[ \widetilde{\psi}(x) - \widetilde{\varphi}(x) \right]\,\mathrm{d}x = 0,$$

and since $\psi_k - \varphi_k \ge 0$, we must have $\widetilde{\psi} - \widetilde{\varphi} \ge 0$.

By the remark following the proof of the bounded convergence theorem, we conclude that $\widetilde{\psi} - \widetilde{\varphi} = 0$ a.e., and therefore $\widetilde{\psi} = \widetilde{\varphi} = f$ a.e., which proves that $f$ is measurable. Finally, since $\varphi_k \to f$ almost everywhere, we have (by definition)

$$\lim_{k \to \infty} \int_{[a,b]}^{\mathcal{L}} \varphi_k(x)\,\mathrm{d}x = \int_{[a,b]}^{\mathcal{L}} f(x)\,\mathrm{d}x,$$

and by $(***)$ and equations (1.9) we see that $\int_{[a,b]}^{\mathcal{R}} f(x)\,\mathrm{d}x = \int_{[a,b]}^{\mathcal{L}} f(x)\,\mathrm{d}x$ as desired.   $\square$

### 1.10.3   Stage Three: Nonnegative Functions

We proceed with the integrals of functions that are measurable and non-negative but not necessarily bounded. It will be important to allow these functions to be extended-valued, that is, these functions may take on the value $+\infty$ (on a measurable set). We recall in this connection the convention that one defines the supremum of a set of positive numbers to be $+\infty$ if the set is unbounded.

In the case of such a function $f$ we define its (extended) *Lebesgue integral* by

$$\int f(x)\,dx = \sup_g \int g(x)\,dx,$$

where the supremum is taken over all measurable functions $g$ such that $0 \le g \le f$, and where $g$ is bounded and supported on a set of finite measure.

**Remark:** With the above definition of the integral, there are only two possible cases: the supremum is either finite, or infinite. In the first case, when $\int f(x)dx < \infty$, we shall say that $f$ is *Lebesgue integrable* (or simply *integrable*). Clearly, if $E$ is any measurable subset of $\mathbb{R}^d$, and $f \ge 0$, then $f\chi_E$ is also positive, and we define

$$\int_E f(x)\,dx = \int f(x)\chi_E(x)\,dx.$$

**Example 29.** *Let*

$$f(x) = \begin{cases} |x|^{-\alpha} & \text{if } |x| \le 1, \\ 0 & \text{otherwise,} \end{cases} \qquad \text{and} \qquad g(x) = \begin{cases} |x|^{-\alpha} & \text{if } |x| > 1, \\ 0 & \text{otherwise.} \end{cases}$$

*Then $f$ is integrable on $\mathbb{R}^d$ if and only if $\alpha < d$ and $g$ is integrable on $\mathbb{R}^d$ if and only if $\alpha > d$.* ✾

The following notation is pretty standard and we will use it in some of the results that follow: $f_n \nearrow f$ refers to a sequence $\{f_n\}$ of monotonically increasing functions that are converging to the limit $f$ as $n \to \infty$ a.e. $x$. We denote the decreasing analogues by $f_n \searrow f$.

**Proposition 21.** *The integral of nonnegative measurable functions enjoys the following properties:*

i) (LINEARITY) *If $f, g \geq 0$ and $\alpha, \beta \in \mathbb{R}$, then*

$$\int (\alpha f + \beta g) = \alpha \int f + \beta \int g.$$

ii) (ADDITIVITY) *If $E$ and $F$ are disjoint subsets of $\mathbb{R}^d$, and $f \geq 0$, then*

$$\int_{E \cup F} f = \int_E f + \int_F f.$$

iii) (MONOTONICITY) *If $0 \leq f \leq g$, then*

$$\int f \leq \int g.$$

iv) *If $g$ is integrable and $0 \leq f \leq g$, then $f$ is integrable.*

v) *If $f$ is integrable, then $f(x) < \infty$ for almost every $x$.*

vi) *If $\int f = 0$, then $f(x) = 0$ for almost every $x$.*

*Proof.* Of the first four assertions, only i) is not an immediate consequence of the definitions, and to prove it we argue as follows: We take $a = b = 1$ and note that if $\varphi \leq f$ and $\psi \leq g$, where both $\varphi$ and $\psi$ are bounded and supported on sets of finite measure, then $\varphi + \psi \leq f + g$, and $\varphi + \psi$ is also bounded and supported on a set of finite measure. Consequently

$$\int f + \int g \leq \int (f + g).$$

To prove the reverse inequality, suppose $\eta$ is bounded and supported on a set of finite measure, and $\eta \leq f + g$. If we define

$$\eta_1(x) = \min\{f(x), \eta(x)\} \qquad \text{and} \qquad \eta_2 = \eta - \eta_1,$$

we note that

$$\eta_1 \leq f \qquad \text{and} \qquad \eta_2 \leq g.$$

Moreover both $\eta_1, \eta_2$ are bounded and supported on sets of finite measure. Hence,

$$\int \eta = \int (\eta_1 + \eta_2) = \int \eta_1 + \int \eta_2 \leq \int f + \int g.$$

Taking the supremum over $\eta$ yields the required inequality.

To prove the conclusion v) we argue as follows: Suppose

$$E_k = \{x \mid f(x) \geq k\} \qquad \text{and} \qquad E_\infty = \{x \mid f(x) = \infty\}.$$

Then,

$$\int f \geq \int \chi_{E_k} f \geq k\, m(E_k),$$

so that $m(E_k) \to 0$ as $k \to \infty$. Since $E_k \searrow E_\infty$, by a previous corollary we have that $m(E_\infty) = 0$. $\qquad\qquad\square$

We now turn our attention to some important convergence theorems for the class of nonnegative measurable functions. To motivate the results that follow, we ask the following question: Suppose $f_n \geq 0$ and $f_n(x) \to f(x)$ for almost every $x$. Is it true that $\int f_n\, dx \to \int f\, dx$?

Unfortunately, the example that follows provides a negative answer to this, and shows that we must change our formulation of the question to obtain a positive convergence result:

**Example 30.** *Let*

$$f_n(x) = \begin{cases} n & \text{if} \quad 0 < x < 1/n, \\ 0 & \text{otherwise.} \end{cases}$$

*Then $f_n(x) \to 0$ for all $x$, yet $\int f_n(x)\, dx = 1$ for all $n$.*

In this particular example, the limit of the integrals is greater than the integral of the limit function. This turns out to be the case in general, as we shall see now:

**Lemma 12 ( Fatou's Lemma).** *Suppose $\{f_n\}$ is a sequence of measurable functions with $f_n \geq 0$. If $\lim_{n \to \infty} f_n(x) = f(x)$ for a.e. $x$, then $\int f \leq \liminf_{n \to \infty} \int f_n$.*

*Proof.* Suppose $0 \leq g \leq f$, where $g$ is bounded and supported on a set $E$ of finite measure. If we set $g_n(x) = \min\{g(x), f_n(x)\}$, then $g_n$ is measurable, supported on $E$, and $g_n(x) \to g(x)$ a.e., so that by the bounded convergence theorem,

$$\int g_n \to \int g.$$

By construction, we also have $g_n \leq f_n$, so that $\int g_n \leq \int f_n$ by the monotonicity of the integral. Thus,

$$\int g \leq \liminf_{n \to \infty} \int f_n.$$

Taking the supremum over all $g$ yields the desired inequality. $\qquad\square$

**Remark:** Note that in the results just presented we do not exclude the cases $\int f = \infty$, or $\liminf_{n \to \infty} f_n = \infty$.

We can now immediately deduce the following series of corollaries:

**Corollary 20.** *Suppose $f$ is a nonnegative measurable function, and $\{f_n\}$ a sequence of nonnegative measurable functions with $f_n(x) \leq f(x)$ and $f_n(x) \to f(x)$ for almost every $x$. Then*

$$\lim_{n \to \infty} \int f_n = \int f.$$

*Proof.* Since $f_n(x) \leq f(x)$ a.e. $x$, we necessarily have $\int f_n \leq \int f$ for all $n$. Hence,

$$\liminf_{n \to \infty} \int f_n \leq \int f.$$

This inequality combined with Fatou's lemma proves the desired limit. $\qquad\square$

In particular, we can now obtain a basic convergence theorem for the class of nonnegative measurable functions:

**Corollary 21** (**Monotone Convergence Theorem**). *Suppose $\{f_n\}$ is a sequence of non-negative measurable functions with $f_n \nearrow f$. Then*

$$\lim_{n \to \infty} \int f_n = \int f.$$

The proof follows immediately from the preceding corollary and its proof. Now, the monotone convergence theorem has the following useful consequence:

**Corollary 22.** *Consider a series $\sum_{k=1}^{\infty} a_k(x)$, where $a_k \geq 0$ is measurable for every $k \geq 1$. Then*

$$\int \sum_{k=1}^{\infty} a_k(x)\, dx = \sum_{k=1}^{\infty} \int a_k(x)\, dx.$$

*If $\sum_{k=1}^{\infty} \int a_k(x)\, dx$ is finite, then the series $\sum_{k=1}^{\infty} a_k(x)$ converges for a.e. $x$.*

*Proof.* Let $f_n(x) = \sum_{k=1}^{n} a_k(x)$ and $f(x) = \sum_{k=1}^{\infty} a_k(x)$. The functions $f_n$ are measurable, $f_n(x) \leq f_{n+1}(x)$, and $f_n(x) \to f(x)$ as $n$ tends to infinity. Since

$$\int f_n = \sum_{k=1}^{n} \int a_k(x)\, dx,$$

the monotone convergence theorem implies

$$\sum_{k=1}^{\infty} \int a_k(x)\, dx = \int \sum_{k=1}^{\infty} a_k(x)\, dx.$$

If $\sum \int a_k < \infty$, then the above implies that $\sum_{k=1}^{\infty} a_k(x)$ is integrable, and by our earlier observation, we conclude that $\sum_{k=1}^{\infty} a_k(x)$ is finite almost everywhere. $\qquad\square$

**Lemma 13 (Borel-Cantelli Lemma).** *If $E_1, E_2, \ldots$ is a collection of measurable subsets with $\sum m(E_k) < \infty$, then the set of points that belong to infinitely many sets $E_k$ has measure zero.*

*Proof.* Let $a_k(x) = \chi_{E_k}(x)$, and note that a point $x$ belongs to infinitely many sets $E_k$ if and only if $\sum_{k=1}^{\infty} a_k(x) = \infty$. Our assumption on $\sum m(E_k)$ says precisely that $\sum_{k=1}^{\infty} \int a_k(x)\, dx < \infty$, and the preceding corollary implies that $\sum_{k=1}^{\infty} a_k(x)$ is finite except possibly on a set of measure zero. $\qquad\square$

### 1.10.4   Stage Four: General Case

If $f$ is any real-valued measurable function on $\mathbb{R}^d$, we say that $f$ is Lebesgue integrable (or just integrable) if the nonnegative measurable function $|f|$ is integrable in the sense

discussed on stage three. If $f$ is Lebesgue integrable, we give a meaning to its integral as follows: First, we may define

$$f^+(x) = \max\{f(x), 0\} \qquad \text{and} \qquad f^-(x) = \max\{-f(x), 0\},$$

so that both $f^+$ and $f^-$ are non-negative and

$$f^+ - f^- = f.$$

Since $f^\pm \leq |f|$, both functions $f^+$ and $f^-$ are integrable whenever $f$ is, and we then define the **Lebesgue integral** of $f$ by

$$\int f = \int f^+ - \int f^-.$$

In practice one encounters many decompositions $f = f_1 - f_2$, where $f_1, f_2$ are both nonnegative integrable functions, and one would expect that regardless of the decomposition of $f$, we always have

$$\int f = \int f_1 - \int f_2.$$

In other words, the definition of the integral should be independent of the decomposition $f = f_1 - f_2$.

Simple applications of the definition and the properties shown previously yield all the elementary properties of the integral:

**Proposition 22.** *The integral of Lebesgue integrable functions is linear, additive, monotonic, and satisfies the triangle inequality.*

**Proposition 23.** *Suppose $f$ is integrable on $\mathbb{R}^d$. Then for every $\varepsilon > 0$:*

    *i) There exists a set of finite measure $B$ (a ball, for example) such that $\int_{B^c} |f| < \varepsilon$.*

    *ii) There is a $\delta > 0$ such that $\int_E |f| < \varepsilon$ whenever $m(E) < \delta$.*

*Condition ii) is known as **absolute continuity**.*

*Proof.* By replacing $f$ with $|f|$ we may assume WLOG that $f \geq 0$. To prove condition i), let $B_N$ denote the ball of radius $N$ centered at the origin, and note that if $f_N(x) = f(x)\chi_{B_N}(x)$, then $f_N \geq 0$ is measurable, $f_N(x) \leq f_{N+1}(x)$, and $\lim_{N \to \infty} f_N(x) = f(x)$. By the monotone convergence theorem, we must have

$$\lim_{N \to \infty} \int f_N = \int f.$$

In particular, for some large $N$,

$$0 \leq \int f - \int f\chi_{B_N} < \varepsilon,$$

and since $1 - \chi_{B_N} = \chi_{B_N^c}$, this implies $\int_{B_N^c} f < \varepsilon$, as we set out to prove.

Now to prove condition ii), we assume again that $f \geq 0$ and we let $f_N(x) = f(x)\chi_{E_N}$, where $E_N = \{x \mid f(x) \leq N\}$. Once again, $f_N \geq 0$ is measurable, $f_N(x) \leq f_{N+1}(x)$, and given $\varepsilon > 0$ there exists (by the monotone convergence theorem) an integer $N > 0$ such that

$$\int (f - f_N) < \frac{\varepsilon}{2}.$$

We now pick $\delta > 0$ so that $N\delta < \varepsilon/2$. If $m(E) < \delta$, then

$$
\begin{aligned}
\int_E f &= \int_E (f - f_N) + \int_E f_N \\
&\leq \int (f - f_N) + \int_E f_N \\
&\leq \int (f - f_N) + N\,m(E) \\
&\leq \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon. \qquad \square
\end{aligned}
$$

Intuitively, integrable functions should in some sense vanish at infinity since their integrals are finite; condition i) of the above proposition attaches a precise meaning to this intuition. One should observe, however, that integrability need not guarantee the more naive pointwise vanishing as $|x|$ becomes large (see [Stein, 2005, Chpt 2, Ex 6]).

We are now ready to prove a cornerstone of the theory of Lebesgue integration, the *Dominated Convergence Theorem*. It can be viewed as a culmination of our efforts, and is a general statement about the interplay between limits and integrals.

**Theorem 63 (Dominated Convergence Theorem).** *Suppose $\{f_n\}$ is a sequence of measurable functions such that $f_n(x) \to f(x)$ a.e. $x$ as $n \to \infty$. If $|f_n(x)| \le g(x)$, where $g$ is integrable, then*

$$\int |f_n - f| \to 0 \qquad \text{as } n \to \infty,$$

*and consequently,*

$$\int f_n \to \int f \qquad \text{as } n \to \infty.$$

*Proof.* For each $N \ge 0$ let $E_N = \{x \mid |x| \le N, g(x) \le N\}$. Given $\varepsilon > 0$, we may argue (as in the first part of above proposition) that there exists $N$ so that

$$\int_{E_N^c} |f| < \varepsilon.$$

Then the functions $f_n \chi_{E_N}$ are bounded (by $N$) and supported on a set of finite measure, so that by the bounded convergence theorem, we have

$$\int_{E_N} |f_n - f| < \varepsilon \qquad \text{for all large } n.$$

Hence, we obtain the estimate

$$\int |f_n - f| = \int_{E_N} |f_n - f| + \int_{E_N^c} |f_n - f|$$

$$\le \int_{E_N} |f_n - f| + 2 \int_{E_N^c} g$$

$$\le \varepsilon + 2\varepsilon = 3\varepsilon$$

for all large $n$. This proves the theorem. $\qquad\square$

## 1.11   Differentiation

### 1.11.1   The Averaging Problem

The following question is referred to as the *averaging problem*:

Suppose $f$ is integrable on $\mathbb{R}^d$. Is it true that

$$\lim_{\substack{m(B) \to 0 \\ x \in B}} \frac{1}{m(B)} \int_B f(y)\, dy = f(x), \quad \text{for a.e. } x\,?$$

Note of course that in the special case when $f$ is continuous at $x$, the limit does converge to $f(x)$. Indeed, given $\varepsilon > 0$, there exists $\delta > 0$ such that $\|f(x) - f(y)\| < \varepsilon$ whenever $\|x - y\| < \delta$.

Since

$$f(x) - \frac{1}{m(B)} \int_B f(y)\, dy = \frac{1}{m(B)} \int_B (f(x) - f(y))\, dy,$$

we find that whenever $B$ is a ball of radius $< \delta/2$ that contains $x$, then

$$\left\| f(x) - \frac{1}{m(B)} \int_B f(y)\, dy \right\| \leq \frac{1}{m(B)} \int_B \|f(x) - f(y)\|\, dy < \varepsilon.$$

as desired.

The averaging problem has an affirmative answer, but to establish that fact, which is qualitative in nature, we need to make some quantitative estimates bearing on the overall behavior of the averages of $f$. This will be done in terms of the maximal averages of $\|f\|$, to which we now turn.

### 1.11.2   Hardy-Littlewood Maximal Function

If $f$ is integrable on $\mathbb{R}^d$, we define its *maximal function* $f^*$ by

$$f^*(x) = \sup_{B:\, x \in B} \frac{1}{m(B)} \int_B \|f(y)\|\, dy, \qquad x \in \mathbb{R}^d,$$

where the supremum is taken over all balls containing the point $x$. In other words, we replace the limit in the statement of the averaging problem by a supremum, and $f$ by its absolute value.

The main properties of $f^*$ we shall need are summarized in the following theorem:

**Theorem 64.** *Suppose $f$ is integrable on $\mathbb{R}^d$. Then,*

   *i)* $f^*$ *is measurable.*

   *ii)* $f^*(x) < \infty$ *for a.e. $x$.*

   *iii)* $f^*$ *satisfies*

$$m(\{x \in \mathbb{R}^d \mid f^*(x) > \alpha\}) \leq \frac{3^d}{\alpha} \|f\|_{L^1(\mathbb{R}^d)} \quad \forall\, \alpha > 0. \tag{1.10}$$

*Proof of i).* The assertion that $f^*$ is measurable is pretty simple. Indeed, the set $E_\alpha = \{x \in \mathbb{R}^d : f^*(x) > \alpha\}$ is open, because if $\overline{x} \in E_\alpha$, there exists a ball $B$ such that $\overline{x} \in B$ and

$$\frac{1}{m(B)} \int_B \|f(y)\|\ \mathrm{d}y > \alpha.$$

Now any point $x$ close enough to $\overline{x}$ will also belong to $B$; hence $x \in E_\alpha$ as well. $\qquad\square$

*Proof of ii).* This condition follows directly from *iii*) once we observe that

$$\{x \mid f^*(x) = \infty\} \subset \{x \mid f^*(x) > \alpha\} \qquad \forall\, \alpha.$$

Taking the limit as $\alpha$ tends to infinity, the third property yields $m(\{x \mid f^*(x) = \infty\}) = 0.$ $\square$

*Proof of iii).* The proof of inequality (1.10) relies on an elementary version of a Vitali covering argument, which is stated in the following lemma aside:

Aside Lemma

**Lemma 14.** *Suppose $B = \{B_1, B_2, \ldots, B_N\}$ is a finite collection of open balls in $\mathbb{R}^d$. Then there exists a disjoint sub-collection $B_{i_1}, B_{i_2}, \ldots, B_{i_k}$ of $B$ that satisfies*

$$m\left(\bigcup_{l=1}^{N} B_l\right) \leq 3^d \sum_{j=1}^{k} m(B_{i_j}). \tag{1.11}$$

Loosely speaking, the lemma tells us that we may always find a disjoint subcollection of balls that covers a fraction of the region covered by the original collection of balls.

Now the proof of iii) is within reach. If we let $E_\alpha = \{x \colon f^*(x) > \alpha\}$, then for each $x \in E_\alpha$ there exists a ball $B_x$ that contains $x$, and such that

$$\frac{1}{m(B_x)} \int_{B_x} \|f(y)\| \, dy > \alpha.$$

Therefore, for each ball $B_x$ we have

$$m(B_x) < \frac{1}{\alpha} \int_{B_x} \|f(y)\| \, dy. \tag{1.12}$$

Fix a compact subset $K$ of $E_\alpha$. Since $K$ is covered by $\cup_{x \in E_\alpha} B_x$, we may select a finite subcover of $K$, say $K \subset \cup_{l=1}^{N} B_l$. The covering lemma discussed above guarantees the existence of a subcollection $B_{i_1}, B_{i_2}, \ldots, B_{i_k}$ of disjoint balls that satisfies inequality (1.11). Now since these balls $B_{i_1}, B_{i_2}, \ldots, B_{i_k}$ are disjoint and satisfy (1.11) as well as (1.12), we find that

$$m(K) \leq m\left(\bigcup_{l=1}^{N} B_l\right) \leq 3^d \sum_{j=1}^{k} m(B_{i_j}) \leq \frac{3^d}{\alpha} \sum_{j=1}^{k} \int_{B_{i_j}} \|f(y)\| \, dy$$

$$= \frac{3^d}{\alpha} \int_{\cup_{j=1}^{k} B_{i_j}} \|f(y)\| \, dy$$

$$\leq \frac{3^d}{\alpha} \int_{\mathbb{R}^d} \|f(y)\| \, dy.$$

Since this inequality is true for all compact subsets $K$ of $E_\alpha$, the proof of the weak type inequality for the maximal operator is complete. $\qquad \square$

**Remark:** Let us clarify the nature of the main conclusion iii). As we shall observe later when we prove the *Lebesgue Differentiation Theorem*, one has that $f^*(x) \geq \|f(x)\|$ for a.e. $x$; the effect of iii) is that, broadly speaking, $f^*$ is not much larger than $\|f\|$. From this point of view, we would have liked to conclude that $f^*$ is integrable, as a result of the assumed integrability of $f$. However, this is not the case, and *iii)* is the best substitute available (we will see more on this on the problems at the end of the chapter).

An inequality of the type (1.10) is called a ***weak-type inequality*** because it is weaker than the corresponding inequality for the $L^1$-norms. Indeed, this can be seen from the ***Tchebychev inequality***, which states that for an arbitrary integrable function $g$,

$$m(\{x\colon \|g(x)\| \geq \alpha\}) \leq \frac{1}{\alpha}\|g\|_{L^1(\mathbb{R}^d)} \quad \forall\, \alpha > 0.$$

We should add that the value of $3^d$ in the inequality (1.10) is unimportant for us. What matters is that this constant be independent of $\alpha$ and $f$.

The estimate obtained for the maximal function now leads to a solution of the averaging problem:

**Theorem 65 (The Lebesgue Differentiation Theorem).** *If $f$ is integrable on $\mathbb{R}^d$, then*

$$\lim_{\substack{m(B)\to 0 \\ x\in B}} \frac{1}{m(B)} \int_B f(y)\,\mathrm{d}y = f(x), \quad \text{for a.e. } x. \tag{1.13}$$

*Proof.* It suffices to show that for each $\alpha > 0$ the set

$$E_\alpha = \left\{ x : \limsup_{\substack{m(B)\to 0 \\ B\,|\,x\in B}} \left\| \frac{1}{m(B)} \int_B f(y)\,\mathrm{d}y - f(x) \right\| > 2\alpha \right\}$$

has measure zero, because this assertion then guarantees that the set $E = \cup_{n=1}^\infty E_{1/n}$ has measure zero, and the limit in (1.13) holds at all points of $E^c$. We fix $\alpha$, and invoke a previous theorem, which states that for each $\varepsilon > 0$ we may select a continuous function $g$ of compact support with $\|f - g\|_{L^1(\mathbb{R}^d)} < \varepsilon$. Then since $g$ is continuous, (1.13) holds not just for a.e. $x$, but in fact for all $x$ (this is shown on our earlier discussion of the averaging problem). That is,

$$\lim_{\substack{m(B)\to 0 \\ x\in B}} \frac{1}{m(B)} \int_B g(y)\,\mathrm{d}y = g(x), \quad \text{for all } x.$$

Since we may write the difference $\frac{1}{m(B)} \int_B f(y)\,dy - f(x)$ as

$$\frac{1}{m(B)} \int_B (f(y) - g(y))\,dy + \frac{1}{m(B)} \int_B g(y)\,dy - g(x) + g(x) - f(x),$$

we find that

$$\limsup_{\substack{m(B) \to 0 \\ B \mid x \in B}} \left\| \frac{1}{m(B)} \int_B f(y)\,dy - f(x) \right\| \leq (f - g)^*(x) + \|g(x) - f(x)\|,$$

where the symbol $*$ indicates the maximal function we previously defined. Consequently, if

$$F_\alpha = \{x : (f - g)^*(x) \geq \alpha\} \quad \text{and} \quad G_\alpha = \{x : \|f(x) - g(x)\| \geq \alpha\},$$

then $E_\alpha \subset (F_\alpha \cup G_\alpha)$, because if $u_1$ and $u_2$ are positive, then $u_1 + u_2 > 2\alpha$ only if $u_i > \alpha$ for at least one $u_i$. On the one hand, *Tchebychev's Inequality* yields

$$m(G_\alpha) \leq \frac{1}{\alpha} \|f - g\|_{L^1(\mathbb{R}^d)},$$

and on the other hand, the weak type estimate for the maximal function gives

$$m(F_\alpha) \leq \frac{3^d}{\alpha} \|f - g\|_{L^1(\mathbb{R}^d)}.$$

The function $g$ was selected so that $\|f - g\|_{L^1(\mathbb{R}^d)} < \varepsilon$. Hence we get

$$m(E_\alpha) \leq \frac{3^d}{\alpha}\varepsilon + \frac{1}{\alpha}\varepsilon.$$

Since $\varepsilon$ is arbitrary, we must have $m(E_\alpha) = 0$, and the proof of the theorem is complete. $\quad\square$

**Remark:** Note that as an immediate consequence of the theorem applied to $\|f\|$, we see that $f^*(x) \geq \|f(x)\|$ for a.e. $x$.

We have worked so far under the assumption that $f$ is integrable. This "global" assumption is slightly out of place in the context of a "local" notion like differentiability. Indeed, the limit in Lebesgue's theorem is taken over balls that shrink to the point $x$, so the behavior of $f$ far from $x$ is irrelevant. Thus, we expect the result to remain valid if we simply assume integrability of $f$ on every ball, as we shall see next.

**Definition 51.** *A measurable function $f$ on $\mathbb{R}^d$ is said to be **locally integrable** if for every ball $B$, the function $f(x)\chi_B(x)$ is integrable. We shall denote the space of all locally integrable functions by $L^1_{\text{loc}}(\mathbb{R}^d)$.*

Loosely speaking, the behavior at infinity does not affect the local integrability of a function. For example, the functions $e^{\|x\|}$ and $\|x\|^{-1/2}$ are both locally integrable, but not integrable on $\mathbb{R}^d$.

Now that we have the notion of local integrability, we can see that the *Lebesgue Differentiation Theorem* holds under weaker assumptions:

**Theorem 66.** *If $f \in L^1_{\text{loc}}(\mathbb{R}^d)$, then*

$$\lim_{\substack{m(B)\to 0 \\ x\in B}} \frac{1}{m(B)} \int_B f(y)\,\mathrm{d}y = f(x), \quad \text{for a.e. } x.$$

**Definition 52.** *If $E$ is a measurable set and $x \in \mathbb{R}^d$, we say that $x$ is a point of **Lebesgue density** of $E$ if*

$$\lim_{\substack{m(B)\to 0 \\ x\in B}} \frac{m(B \cap E)}{m(B)} = 1$$

*Loosely speaking, this condition says that small balls around $x$ are almost entirely covered by $E$. More precisely, for every $\alpha < 1$ close to 1, and every ball of sufficiently small radius containing $x$, we have*

$$m(B \cap E) \geq \alpha m(B).$$

*Thus $E$ covers at least a proportion $\alpha$ of $B$.*

**Corollary 23.** *Suppose $E$ is a measurable subset of $\mathbb{R}^d$. Then:*

   *i) Almost every $x \in E$ is a point of density of $E$.*

   *ii) Almost every $x \notin E$ is not a point of density of $E$.*

**Definition 53.** *If $f$ is locally integrable on $\mathbb{R}^d$, the **Lebesgue set** of $f$ consists of all points $\hat{x} \in \mathbb{R}^d$ for which $f(\hat{x})$ is finite and*

$$\lim_{\substack{m(B)\to 0 \\ \hat{x}\in B}} \frac{1}{m(B)} \int_B \|f(y) - f(\hat{x})\| \, dy = 0.$$

**Remark:** At this stage, two simple observations about this definition are in order. First, $\hat{x}$ belongs to the Lebesgue set of $f$ whenever $f$ is continuous at $\hat{x}$. Second, if $\hat{x}$ is in the Lebesgue set of $f$, then

$$\lim_{\substack{m(B)\to 0 \\ \hat{x}\in B}} \frac{1}{m(B)} \int_B f(y) \, dy = f(\hat{x}).$$

**Corollary 24.** *If $f$ is locally integrable on $\mathbb{R}^d$, then almost every point belongs to the Lebesgue set of $f$.*

### 1.11.3  Functions of Bounded Variation

**Definition 54.** *A function $F$ is said to be of **bounded variation** on an interval $[a, b]$ if the variations of $F$ over all partitions of such interval are bounded; that is, there exists $M < \infty$ so that*

$$\sum_{j=1}^{N} \|F(t_j) - F(t_{j-1})\| \leq M$$

*for all partitions $a = t_0 < t_1 < \cdots < t_N = b$.*

Notice that if $F$ is real-valued, monotonic, and bounded, then $F$ is of bounded variation. Indeed, if for example $F$ is nondecreasing and bounded by $M$, we see that

$$\sum_{j=1}^{N} \|F(t_j) - F(t_{j-1})\| = \sum_{j=1}^{N} F(t_j) - F(t_{j-1})$$
$$= F(b) - F(a) \leq 2M.$$

Notice also that if $F$ is differentiable at every point, and $F'$ is bounded, then $F$ is of bounded variation. Indeed, if $\|F'\| \leq M$, then the mean value theorem implies

$$\|F(x) - F(y)\| \leq M \|x - y\|, \qquad \forall\, x, y \in [a, b],$$

which in turn implies that $\sum_{j=1}^{N} \|F(t_j) - F(t_{j-1})\| \leq M(b - a)$.

**Example 31.** *Look at the function*

$$F(x) = \begin{cases} x^a \sin(x^{-b}) & \text{if } 0 < x \leq 1 \\ 0 & \text{if } x = 0. \end{cases}$$

*We have that $F$ is of bounded of variation if and only if $a > b$.*

**Definition 55.** *The **total variation** of a function $F$ on $[a, x]$ (where $a \leq x \leq b$) is defined by*

$$T_F(a, x) = \sup \sum_{j=1}^{N} \|F(t_j) - F(t_{j-1})\|,$$

*where the supremum is taken over all partitions of $[a, x]$.*

**Definition 56.** *The **positive variation** of a function $F$ on $[a, x]$ (where $a \leq x \leq b$) is defined by*

$$P_F(a, x) = \sup \sum_{(+)} (F(t_j) - F(t_{j-1})),$$

*where the sum is taken over all $j$ such that $F(t_j) \geq F(t_{j-1})$, and the sup is taken over all partitions of $[a, x]$.*

**Definition 57.** *The **negative variation** of a function $F$ on $[a, x]$ (where $a \leq x \leq b$) is defined by*

$$N_F(a, x) = \sup \sum_{(-)} (F(t_{j-1}) - F(t_j)),$$

*where the sum is taken over all $j$ such that $F(t_{j-1}) \geq F(t_j)$, and the sup is taken over all partitions of $[a, x]$.*

**Lemma 15.** *Suppose F is real-valued and of bounded variation on $[a, b]$. Then for all $a \leq x \leq b$, we have*

$$F(x) - F(a) = P_F(a, x) - N_F(a, x), \quad \text{and} \quad T_F(a, x) = P_F(a, x) + N_F(a, x).$$

**Theorem 67.** *A real-valued function F on $[a, b]$ is of bounded variation if and only if F is the difference of two increasing bounded functions.*

*Proof.* $(\Leftarrow)$ Clearly, if $F = F_1 - F_2$, where each $F_j$ is bounded and increasing, then $F$ is of bounded variation.

$(\Rightarrow)$ Conversely, suppose $F$ is of bounded variation. Then, we let $F_1(x) = P_F(a, x) + F(a)$ and $F_2(x) = N_F(a, x)$. Clearly, both $F_1$ and $F_2$ are increasing and bounded, and by the above lemma $F(x) = F_1(x) - F_2(x)$. $\qquad \square$

The next result lies at the heart of the theory of differentiation:

**Theorem 68.** *If F is of bounded variation on $[a, b]$, then F is differentiable almost everywhere. In other words, the quotient*

$$\lim_{h \to 0} \frac{F(x + h) - F(x)}{h}$$

*exists for almost every $x \in [a, b]$.*

The proof of this theorem relies on the lemma discussed below as well as on the concept of the so called *Dini* numbers.

**Lemma 16 (Riesz's Lemma).** *Suppose G is real-valued and continuous on $\mathbb{R}$. Let*

$$E = \{x \in \mathbb{R} : G(x + h) > G(x) \text{ for some } h = h_x > 0\}$$

*If E is non-empty, then it must be open, and hence can be written as a countable disjoint union of open intervals $E = \bigcup(a_k, b_k)$. If $(a_k, b_k)$ is a finite interval of this union, then we have that*

$$G(b_k) - G(a_k) = 0.$$

*Proof.* Since $G$ is continuous, it is clear that $E$ is open whenever it is non-empty and can therefore be written as a disjoint union of countably many open intervals. If $(a_k, b_k)$ denotes a finite interval in this decomposition, then $a_k \notin E$; therefore we cannot have $G(b_k) > G(a_k)$.

We now suppose that $G(b_k) < G(a_k)$. By continuity, there exists $a_k < c < b_k$ so that

$$G(c) = \frac{G(a_k) + G(b_k)}{2},$$

and in fact we may choose $c$ farthest to the right in the interval $(a_k, b_k)$. Since $c \in E$, there must exist a $d > c$ in $E$ so that $G(d) > G(c)$. Since $b_k \notin E$, we must have $G(x) \leq G(b_k)$ for all $x \geq b_k$; therefore $d < b_k$. Since $G(d) > G(c)$, there exists (by continuity) $c' > d$ with $c' < b_k$ and $G(c') = G(c)$, which contradicts the fact that $c$ was chosen farthest to the right in $(a_k, b_k)$.

This shows that we must have $G(a_k) = G(b_k)$, and the lemma is proved. $\qquad\square$

**Remark:** This result is usually called the **Rising Sun Lemma** for the following reason. If one thinks of the sun rising from the east (at the right) with the rays of light parallel to the $x$-axis, then the points $(x, G(x))$ on the graph of $G$, with $x \in E$, are precisely the points which are in the shade; these points appear in bold in Figure 1.7.



Figure 1.7: Visual representation of the *Rising Sun Lemma*.

**Corollary 25.** *If $F$ is increasing and continuous, then $F'$ exists almost everywhere. Moreover $F'$ is measurable, non-negative, and*

$$\int_a^b F'(x)\,dx \leq F(b) - F(a).$$

*In particular, if $F$ is bounded on $\mathbb{R}$, then $F'$ is integrable on $\mathbb{R}$.*

**Remark:** Note that if we had equality, then the above corollary would give us the *Fundamental Theorem of Calculus*. However, we cannot go any farther than the inequality above if we allow all continuous increasing functions, as it is illustrated by the following important example, the so-called *Cantor-Lebesgue function*.

### 1.11.4 The Cantor-Lebesgue function

The following simple construction yields a continuous function $F : [0,1] \to [0,1]$ that is increasing with $F(0) = 0$ and $F(1) = 1$, but $F'(x) = 0$ almost everywhere! Hence $F$ is of bounded variation, but

$$\int_a^b F'(x)\, dx \neq F(b) - F(a).$$

Consider the standard triadic Cantor set $\mathcal{C} \subset [0,1]$, where $\mathcal{C} = \bigcap_{k=0}^{\infty} C_k$ and each $C_k$ is a disjoint union of $2^k$ closed intervals. For example, $C_1 = [0,1/3] \cup [2/3,1]$.

Now let $F_1(x)$ be the continuous increasing function on $[0,1]$ (and linear on $C_1$) that satisfies

$$F_1(x) = \begin{cases} 0 & \text{if } x = 0, \\ 1/2 & \text{if } 1/3 \leq x \leq 2/3, \\ 1 & \text{if } x = 1. \end{cases}$$

Similarly, let $F_2(x)$ (see Figure 1.8) be the continuous increasing function on $[0,1]$ (and linear on $C_2$) that satisfies

$$F_2(x) = \begin{cases} 0 & \text{if } x = 0, \\ 1/4 & \text{if } 1/9 \leq x \leq 2/9, \\ 1/2 & \text{if } 1/3 \leq x \leq 2/3, \\ 3/4 & \text{if } 7/9 \leq x \leq 8/9, \\ 1 & \text{if } x = 1. \end{cases}$$

This process yields a sequence of continuous increasing functions $\{F_n\}_{n=1}^{\infty}$ such that clearly

$$\|F_{n+1}(x) - F_n(x)\| \leq \frac{1}{2^{n+1}}.$$

Figure 1.8: Here's a visualization of the construction of $F_2$.

Hence $\{F_n\}_{n=1}^{\infty}$ converges uniformly to a continuous limit $F$, which is called the **Cantor-Lebesgue function** (see Figure 1.9).



Figure 1.9: Here's a visualization of the *Cantor-Lebesgue* function.

By construction $F$ is increasing, $F(0) = 0$, $F(1) = 1$, and we see that $F$ is constant on each interval of the complement of the Cantor set. Since $m(\mathcal{C}) = 0$, we find that $F'(x) = 0$ almost everywhere, as desired.

The considerations in this section, as well as this last example, show that the assumption of bounded variation guarantees the existence of a derivative almost everywhere, but not

the validity of the formula

$$\int_a^b F'(x)\,dx = F(b) - F(a).$$

In order to achieve this equality we need to consider *absolutely continuous* functions, which we turn to next.

## 1.11.5 Absolutely Continuous Functions

**Definition 58.** *A function f defined on some interval $[a, b]$ is said to be* **absolutely continuous** *if for any $\varepsilon > 0$, there exists a $\delta > 0$ such that*

$$\sum_{k=1}^{N} \|f(b_k) - f(a_k)\| < \varepsilon \quad \text{whenever} \quad \sum_{k=1}^{N} (b_k - a_k) < \delta,$$

*where the intervals $(a_k, b_k)$ (for $k = 1, \ldots, N$) are disjoint intervals.*

From the definition, it is clear that absolutely continuous functions are continuous, and in fact uniformly continuous. Also, note that if $F$ is absolutely continuous on a bounded interval, then it is also of bounded variation on the same interval. Moreover, its total variation is continuous (in fact absolutely continuous). As a consequence, the decomposition of such a function $F$ into two monotonic functions given previously in Theorem 67 shows that each of these functions is continuous.

**Remark:** If $F(x) = \int_a^x f(y)\,dy$ where $f$ is integrable, then $F$ is absolutely continuous. This remark shows that absolute continuity is a necessary condition to impose on $F$ if we hope to prove the desired equality $\int_a^b F'(x)\,dx = F(b) - F(a)$.

**Theorem 69.** *If F is absolutely continuous on $[a, b]$, then $F'(x)$ exists almost everywhere. Moreover, if $F'(x) = 0$ for a.e. x, then F is constant.*

The following lemma and corollary are used in the proof of *Theorem 69*. First we need to state the definition of a *Vitali covering*:

**Definition 59.** *A collection $\mathcal{B}$ of balls $\{B\}$ is said to be a **Vitali covering** of a set $E$ if for every $x \in E$ and any $\eta > 0$ there is a ball $B \in \mathcal{B}$, such that $x \in B$ and $m(B) < \eta$. Thus every point is covered by balls of arbitrarily small measure.*

**Lemma 17.** *Suppose $E$ is a set of finite measure and $\mathcal{B}$ is a Vitali covering of $E$. For any $\delta > 0$, we can find finitely many balls $B_1, \ldots, B_N$ in $\mathcal{B}$ that are disjoint and so that*

$$\sum_{i=1}^{N} m(B_i) \geq m(E) - \delta.$$

**Corollary 26.** *We can arrange the choice of the balls from the Vitali covering given in the above lemma so that*

$$m\left( E - \bigcup_{i=1}^{N} B_i \right) < 2\delta.$$

The culmination of all our efforts is contained in the next theorem. In particular, it resolves the problem of establishing the reciprocity between differentiation and integration, known as the *Fundamental Theorem of Calculus*:

**Theorem 70 (Fundamental Theorem of Calculus).** *Suppose $F$ is absolutely continuous on $[a, b]$. Then $F'$ exists almost everywhere and is integrable. Moreover,*

$$F(x) - F(a) = \int_a^x F'(y) \, dy, \qquad \forall\, a \leq x \leq b.$$

*By selecting $x = b$ we get $F(b) - F(a) = \int_a^b F'(y) \, dy$.*

*Conversely, if $f$ is integrable on $[a, b]$, then there exists an absolutely continuous function $F$ such that $F'(x) = f(x)$ almost everywhere, and in fact, we may take $F(x) = \int_a^x f(y) \, dy$.*

*Proof.* Since we know that a real-valued absolutely continuous function is the difference of two continuous increasing functions, Corollary 25 shows that $F'$ is integrable on $[a, b]$. Now let $G(x) = \int_a^x F'(y) \, dy$. Then $G$ is absolutely continuous; hence so is the difference $G(x) - F(x)$. By the *Lebesgue Differentiation Theorem*, we know that $G'(x) = F'(x)$ for a.e. $x$;

hence the difference $F - G$ has derivative $0$ almost everywhere. By Theorem 69, we conclude that $F - G$ is constant, and evaluating this expression at $x = a$ gives the desired result.

The converse is a consequence of the observation we made earlier, namely that $\int_a^x f(y)\,dy$ is absolutely continuous, and the *Lebesgue Differentiation Theorem*, which gives $F(x) = f(x)$ almost everywhere. $\qquad\square$

## 1.12 Hilbert Spaces

**Definition 60.** *A set $\mathcal{H}$ is a **Hilbert space** if it satisfies the following:*

  *i)* *$\mathcal{H}$ is a vector space over $\mathbb{C}$ (or $\mathbb{R}$).*

  *ii)* *$\mathcal{H}$ is equipped with an inner product $\langle \cdot, \cdot \rangle$, so that*

      – *$f \mapsto \langle f, g \rangle$ is linear on $\mathcal{H}$ for every $g \in \mathcal{H}$.*
      – *$\langle f, g \rangle = \overline{\langle g, f \rangle}$.*
      – *$\langle f, f \rangle \geq 0$ for all $f \in \mathcal{H}$.*

  *iii)* *We let $\|f\| = \sqrt{\langle f, f \rangle}$. Then $\|f\| = 0$ if and only if $f = 0$.*

  *iv)* *$\mathcal{H}$ is complete in the metric $d(f, g) = \|f - g\| = \sqrt{\langle f - g, f - g \rangle}$.*

**Remark**: Notice that the Cauchy-Schwarz and triangle inequalities

$$\|\langle f, g \rangle\| \leq \|f\|\|g\| \qquad \text{and} \qquad \|f + g\| \leq \|f\| + \|g\|$$

are in fact easy consequences of assumptions i) and ii) of our definition. Note also that saying that $\mathcal{H}$ a *Hilbert space* is the same as saying that $\mathcal{H}$ is a Banach space (i.e., a complete normed linear space), with the norm induced by an inner product $\langle \cdot, \cdot \rangle$.

**Definition 61.** *If $F$ is a function defined in the unit disc $\mathbb{D}^2$, we say that $F$ has a **radial limit** at the point $-\pi \leq \theta \leq \pi$ on the circle, if the limit*

$$\lim_{\substack{r \to 1 \\ r < 1}} F(re^{i\theta})$$

*exists.*

**Theorem 71.** *A bounded holomorphic function $F(re^{i\theta})$ on the unit disc has radial limits at almost every $\theta$.*

**Definition 62.** *The **Hardy space** $H^2(\mathbb{D}^2)$ is the space that consist of all holomorphic functions $F$ on the unit disc $\mathbb{D}^2$ that satisfy*

$$\sup_{0 \leq r < 1} \sum_{n=0}^{\infty} \|a_n\|^2 r^{2n} = \sup_{0 \leq r < 1} \frac{1}{2\pi} \int_{-\pi}^{\pi} \|F(re^{i\theta})\|^2 \, d\theta < \infty,$$

*where the $a_n$ are the Fourier coefficients for $n \geq 0$ and $o$ for $n < 0$. We also take the norm for functions $F$ in this class, $\|F\|_{H^2(\mathbb{D}^2)}$, to be the square root of the above quantity.*

Note that if $F$ is bounded, then $F \in H^2(\mathbb{D}^2)$, and moreover the conclusion of the existence of radial limits almost everywhere stated in the above theorem holds for any $F \in H^2(\mathbb{D}^2)$. Finally, we note that

$$F \in H^2(\mathbb{D}^2) \qquad \Longleftrightarrow \qquad F(z) = \sum_{n=0}^{\infty} a_n z^n$$

with $\sum_{n=0}^{\infty} \|a_n\|^2 < \infty$. Furthermore, we have

$$\sum_{n=0}^{\infty} \|a_n\|^2 = \|F\|_{H^2(\mathbb{D}^2)}^2.$$

This states in particular that $H^2(\mathbb{D}^2)$ is in fact a Hilbert space that can be viewed as the "subspace" $\ell^2(\mathbb{Z}^+)$ of $\ell^2(\mathbb{Z})$, consisting of all $\{a_n\} \in \ell^2(\mathbb{Z})$, with $a_n = 0$ when $n < 0$.

*Side Note*: The *parallelogram law* states that in a Hilbert space $\mathcal{H}$, we have

$$\|A + B\|^2 + \|A - B\|^2 = 2[\|A\|^2 + \|B\|^2] \qquad \text{for all } A, B \in \mathcal{H}.$$

**Definition 63.** *Let $\mathcal{H}$ and $\mathcal{H}'$ be Hilbert spaces with respective inner products $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ and $\langle \cdot, \cdot \rangle_{\mathcal{H}'}$ and the corresponding norms $\| \cdot \|_{\mathcal{H}}$ and $\| \cdot \|_{\mathcal{H}'}$. A mapping $U \colon \mathcal{H} \to \mathcal{H}'$ between these spaces is called **unitary** if:*

   *i)* *$U$ is linear, i.e. $U(\alpha f + \beta g) = \alpha U(f) + \beta U(g)$.*

   *ii)* *$U$ is a bijection.*

   *iii)* *$\|Uf\|_{\mathcal{H}'} = \|f\|_{\mathcal{H}}$ for all $f \in \mathcal{H}$.*

**Definition 64.** *If $S$ is a subspace of a Hilbert space $\mathcal{H}$, we define the **orthogonal complement** of $S$ by*

$$S^{\perp} = \{ f \in \mathcal{H} \mid \langle f, g \rangle = 0 \quad \forall g \in S \}.$$

Clearly, $S^{\perp}$ is also a subspace of $\mathcal{H}$, and moreover $S \cap S^{\perp} = \{0\}$. To see this, note that if $f \in S \cap S^{\perp}$, then $f$ must be orthogonal to itself; thus $0 = \langle f, f \rangle = \|f\|$, and therefore $f = 0$. Moreover, $S^{\perp}$ is itself a closed subspace. Indeed, if $f_n \to f$, then $\langle f_n, g \rangle \to \langle f, g \rangle$ for every $g$ by the Cauchy-Schwarz inequality. Hence if $\langle f_n, g \rangle = 0$ for all $g \in S$ and all $n$, then $\langle f, g \rangle = 0$ for all those $g$.

**Proposition 24.** *If $S$ is a closed subspace of a Hilbert space $\mathcal{H}$, then*

$$\mathcal{H} = S \oplus S^{\perp}.$$

The notation in the proposition means that every $f \in \mathcal{H}$ can be written uniquely as $f = g + h$, where $g \in S$ and $h \in S^{\perp}$; we then say that $\mathcal{H}$ is the **direct sum** of $S$ and $S^{\perp}$. This is equivalent to saying that any $f \in \mathcal{H}$ is the sum of two elements, one in $S$, the other in $S^{\perp}$, and that $S \cap S^{\perp}$ contains only 0.

With the decomposition $\mathcal{H} = S \oplus S^{\perp}$ one has the natural projection onto $S$ defined by

$$P_s(f) = g, \qquad \text{where} \quad f = g + h \quad \text{and} \quad g \in S, h \in S^{\perp}.$$

The mapping $P_s$ is called the **orthogonal projection** onto $S$ and satisfies the following simple properties:

i) $f \mapsto P_s(f)$ is linear.

ii) $P_s(f) = f$ whenever $f \in S$.

iii) $P_s(f) = 0$ whenever $f \in S^{\perp}$.

iv) $\|P_s(f)\| \leq \|f\|$ for all $f \in \mathcal{H}$.

Property $i)$ means that $P_s(\alpha f_1 + \beta f_2) = \alpha P_s(f_1) + \beta P_s(f_2)$, whenever $f_1, f_2 \in \mathcal{H}$ and $\alpha$ and $\beta$ are scalars.

Now let us look at a very important result in the next example:

**Example 32.** *Consider $L^2([-\pi, \pi])$ and let S denote the subspace that consists of all $F \in L^2([-\pi, \pi])$ with*

$$F(\theta) \sim \sum_{n=0}^{\infty} a_n e^{in\theta}.$$

*In other words, S is the space of square integrable functions whose Fourier coefficients $a_n$ vanish for $n < 0$. From the proof of Fatou's theorem (see Section §1.10), this implies that S can be identified with the Hardy space $H^2(\mathbb{D}^2)$, and so is a closed subspace unitarily isomorphic to $\ell^2(\mathbb{Z}^+)$.*

*Therefore, using this identification, if P denotes the orthogonal projection from $L^2([-\pi, \pi])$ to S, we may also write $P(f)(z)$ for the element corresponding to $H^2(\mathbb{D}^2)$, that is,*

$$P(f)(z) = \sum_{n=0}^{\infty} a_n z^n.$$

*Now given $f \in L^2([-\pi, \pi])$, we define the **Cauchy integral** of f by*

$$C(f)(z) = \frac{1}{2\pi i} \int_{\gamma} \frac{f(\zeta)}{\zeta - z} \, d\zeta,$$

*where $\gamma$ denotes the unit circle and z belongs to the unit disc.*

*Then we have the identity*

$$P(f)(z) = C(f)(z), \qquad \forall z \in \mathbb{D}^2.$$

*Indeed, since $f \in L^2$ it follows by the Cauchy-Schwarz inequality that $f \in L^1([-\pi, \pi])$, and therefore we may interchange the sum and integral in the following calculation (recall $\|z\| < 1$):*

$$
\begin{aligned}
P(f)(z) = \sum_{n=0}^{\infty} a_n z^n &= \sum_{n=0}^{\infty} \left( \frac{1}{2\pi} \int_{-\pi}^{\pi} f(e^{i\theta}) e^{-in\theta} \, d\theta \right) z^n \\
&= \frac{1}{2\pi} \int_{-\pi}^{\pi} f(e^{i\theta}) \sum_{n=0}^{\infty} (e^{-i\theta} z)^n \, d\theta \\
&= \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{f(e^{i\theta})}{1 - e^{-i\theta} z} \, d\theta \\
&= \frac{1}{2\pi i} \int_{-\pi}^{\pi} \frac{f(e^{i\theta})}{e^{i\theta} - z} \, ie^{i\theta} \, d\theta \\
&= C(f)(z).
\end{aligned}
$$

## 1.12.1 Linear Functionals and the Riesz Representation Theorem

It is a remarkable fact that every continuous linear functional on a Hilbert space arises as an inner product, as stated by the following theorem:

**Theorem 72 (Riesz Representation Theorem).** *Let $\ell$ be a continuous linear functional on a Hilbert space $\mathcal{H}$. Then, there exists a unique $g \in \mathcal{H}$ such that*

$$
\ell(f) = \langle f, g \rangle \qquad \forall f \in \mathcal{H}.
$$

*Moreover, $\|\ell\| = \|g\|$.*

*Proof.* Consider the subspace of $\mathcal{H}$ defined by

$$
\mathcal{S} = \{ f \in \mathcal{H} \mid \ell(f) = 0 \}.
$$

Since $\ell$ is continuous the subspace $\mathcal{S}$, which is called the nullspace of $\ell$, is closed. If $\mathcal{S} = \mathcal{H}$, then $\ell = 0$ and we take $g = 0$. Otherwise $\mathcal{S}^{\perp}$ is non-trivial and we may pick any $h \in \mathcal{S}^{\perp}$ with $\|h\| = 1$. With this choice of $h$ we determine $g$ by setting $g = \overline{\ell(h)} h$. Thus if we let $u = \ell(f)h - \ell(h)f$, then $u \in \mathcal{S}$, and therefore $\langle u, h \rangle = 0$. Hence

$$
0 = \langle \ell(f)h - \ell(h)f, h \rangle = \ell(f) \langle h, h \rangle - \langle f, \overline{\ell(h)}h \rangle.
$$

Since $\langle h, h \rangle = 1$, we find that $\ell(f) = \langle f, g \rangle$ as desired. $\qquad \square$

The first application of the Riesz representation theorem is to determine the existence of the "adjoint" of a linear transformation:

**Proposition 25.** *Let $T : \mathcal{H} \to \mathcal{H}$ be a bounded linear transformation. There exists a unique bounded linear transformation $T^*$ on $\mathcal{H}$ so that*

*i)* $\langle Tf, g \rangle = \langle f, T^*g \rangle$.

*ii)* $\|T\| = \|T^*\|$.

*iii)* $(T^*)^* = T$.

*The linear operator $T^* : \mathcal{H} \to \mathcal{H}$ satisfying the above conditions is called the **adjoint** of $T$. In the special case when $T = T^*$ we say that $T$ is **symmetric**.*

# Chapter 2

# Topology

Topology is concerned with those properties that remain invariant under continuous transformations. In the context of Klein's Erlanger Programm (where it receives a brief mention under its old name of *analysis situs*) it is the "geometry" of groups of continuous invertible transformations, or *homeomorphisms*. The "spaces" to which transformations are applied, and indeed the meaning of "continuous," remain somewhat open. When these terms are interpreted in the most general way, as subject only to certain axioms, one has what is known as *general topology* (or *point-set topology*, as it is also known). This is where we will focus most of our attention in these introductory notes, although we will also present the basics of what is known as *algebraic topology* at the very last section of this chapter.

## 2.1 Topological Spaces

**Definition 65.** *If $X$ is a set, a **topology** on $X$ is a collection $\mathcal{T}$ of subsets of $X$ satisfying the following properties:*

*i)* $X$ *and* $\varnothing$ *are elements of* $\mathcal{T}$.

*ii)* $\mathcal{T}$ *is closed under finite intersections: if* $U_1, \ldots, U_n$ *are elements of* $\mathcal{T}$, *then their intersection* $U_1 \cap \cdots \cap U_n$ *is an element of* $\mathcal{T}$ *as well.*

*iii)* $\mathcal{T}$ *is closed under arbitrary unions: if* $(U_\alpha)_{\alpha \in A}$ *is any (finite or infinite) family of elements of* $\mathcal{T}$, *then their union* $\cup_{\alpha \in A} U_\alpha$ *is also an element of* $\mathcal{T}$.

*In other words, a topology on X is a collection of all open sets of X. A pair* $(X, \mathcal{T})$ *consisting of a set X together with a topology* $\mathcal{T}$ *on X is called a **topological space**.*

**Example 33 (Simple Topologies).** *a) Let X be any given set and let* $\mathcal{T}$ *be the collection of all subsets of X. Then* $\mathcal{T}$ *is a topology on X, called the **discrete topology** on X (see Figure 2.1 a)), and* $(X, \mathcal{T})$ *is called a **discrete space**.*



(a) Discrete topology        (b) Trivial topology        (c) $\{\{1\}, \{1, 2\}, \{1, 2, 3\}, \varnothing\}$

Figure 2.1: Topologies on the set $\{1, 2, 3\}$.

*b) Let Y be any set, and let* $\mathcal{T} = \{Y, \varnothing\}$ *(see Figure 2.1 b)). This topology that consists only of the entire space and the empty space is called the **trivial topology** on Y.*

*c) Let Z be the set* $\{1, 2, 3\}$, *and declare the open subsets to be* $\{1\}$, $\{1, 2\}$, $\{1, 2, 3\}$, *and the empty set. This is the topology illustrated on Figure 2.1 c).*

We can easily verify that each of the preceding examples is in fact a topology by checking that they satisfy all three conditions stated on our definition. I'll leave it to you as a fun exercise to do so.

**Remark:** Note that our definition of a topology is based entirely on the concept of open sets, even though we have made no mention whatsoever of a metric function on the space where we define such topology. This is in stark contrast to our previous exposure to open sets in Chapter 1, where everything depended on a metric. In a word, topology is a generalization to the concept of metric spaces, which are indeed simply a special case of topological

spaces, i.e., they are topological spaces whose topologies (collections of open sets) are defined in terms of the metric (balls of some radius). We now make the following definition: a topological space $(X, \mathcal{T})$ is said to be ***metrizable*** if there is a metric $d\colon X \times X \to [0, \infty)$ such that the topology induced by $d$ is $\mathcal{T}$. One of the themes in the study of topology is to find necessary and sufficient conditions for which a topological space can be endowed with a metric. There will be more on this later on; I also urge you to look up (after completing your readings here of course) a famous result known as *Urysohn's Metrization Theorem*. It states that every Hausdorff second countable regular space is metrizable. All these words will start to make sense soon, so don't freak out!

**Definition 66.** *Let $X$ be a set and let*

$$\mathcal{T}_f = \{U \subseteq X \mid X \smallsetminus U \text{ is finite or is all of } X\}.$$

*Then $\mathcal{T}_f$ is a topology on $X$, called the **finite-complement topology**.*

To show that $\mathcal{T}_f$ is indeed a topology, notice that both $X$ and $\varnothing$ are in $\mathcal{T}_f$, since $X \smallsetminus X$ is finite and $X \smallsetminus \varnothing$ is all of $X$. Also, if $\{U_\alpha\}$ is an indexed family of nonempty elements of $\mathcal{T}_f$, we have that $\cup U_\alpha \subset \mathcal{T}_f$. To see why, compute

$$X \smallsetminus \bigcup U_\alpha = \bigcap (X \smallsetminus U_\alpha).$$

The latter set is finite because each set $X \smallsetminus U_\alpha$ is finite. Finally, if $U_1, \ldots, U_n$ are nonempty elements of $\mathcal{T}_f$, to show that $\cap U_i \subset \mathcal{T}_f$, we compute

$$X \smallsetminus \bigcap_{i=1}^{n} U_i = \bigcup_{i=1}^{n} (X \smallsetminus U_i).$$

The latter set is a finite union of finite sets, therefore it is finite.

**Definition 67.** *Let $X$ be a set and let*

$$\mathcal{T}_C = \{U \subseteq X \mid X \smallsetminus U \text{ is countable or is all of } X\}.$$

*Then $\mathcal{T}_C$ is also a topology on $X$, called the **cocountable topology** (or also the **countable-complement topology**).*

Suppose that $X$ is a topological space and $A$ is any subset of $X$. We now define several related subsets that we previously discussed on Chapter 1, although now we are presenting them in the more general setting of topology:

- The ***closure*** of $A$ in $X$, denoted by $\bar{A}$, is the set

$$\bar{A} = \bigcap \{B \subseteq X \mid A \subseteq B, B \text{ is closed in } X\}.$$

- The ***interior*** of $A$, denoted by $\text{Int}(A)$ or $\mathring{A}$, is given by

$$\mathring{A} = \bigcup \{C \subseteq X \mid C \subseteq A, C \text{ is open in } X\}.$$

- The ***exterior*** of $A$, denoted by $\text{Ext}(A)$, is given by $\text{Ext}(A) = X \smallsetminus \bar{A}$.

- The ***boundary*** of $A$, denoted by $\partial A$, is given by $\partial A = X \smallsetminus \text{Int}(A) \cup \text{Ext}(A))$.



Figure 2.2: Interior, exterior, and boundary points.

**Remark:** It follows immediately from the properties of open and closed subsets that $\bar{A}$ is closed and $\mathring{A}$ is open. As we previously mentioned in Chapter 1, $\bar{A}$ is the smallest closed subset containing $A$ while $\mathring{A}$ is the largest open subset contained in $A$. Notice also from the above definitions that for any subset $A \subseteq X$, the whole space $X$ is equal to the disjoint union of $\mathring{A}$, $\text{Ext}(A)$, and $\partial A$. Moreover, $\mathring{A}$ and $\text{Ext}(A)$ are open in $X$, while $\bar{A}$ and $\partial A$ are closed in $X$.

It is also sometimes useful to compare different topologies on the same set:

**Definition 68.** *Given two topologies $\mathcal{T}_1$ and $\mathcal{T}_2$ on a set $X$, we say that $\mathcal{T}_1$ is **finer** than $\mathcal{T}_2$ if $\mathcal{T}_2 \subseteq \mathcal{T}_1$, and **coarser** than $\mathcal{T}_2$ if $\mathcal{T}_1 \subseteq \mathcal{T}_2$. We say that $\mathcal{T}_1$ is **comparable** to $\mathcal{T}_2$ if either $\mathcal{T}_1 \subseteq \mathcal{T}_2$ or $\mathcal{T}_2 \subseteq \mathcal{T}_1$.*

The terminology in this definition is meant to suggest the picture of a subset that is open in a coarser topology being further subdivided into smaller open subsets in a finer topology. It can be shown that the identity map of $X$ is continuous as a map from $(X, \mathcal{T}_1)$ to $(X, \mathcal{T}_2)$ if and only if $\mathcal{T}_1$ is finer than $\mathcal{T}_2$, and furthermore it is a homeomorphism if and only if $\mathcal{T}_1 = \mathcal{T}_2$.

Here are a few explicit examples of homeomorphisms that we should keep in mind:

**Example 34.** *a) Any open ball in $\mathbb{R}^n$ is homeomorphic to any other open ball: the homeomorphism can easily be constructed as a composition of translations $x \mapsto x + x_0$ and dilations $x \mapsto \alpha x$ (for some scalar $\alpha$). Similarly, all spheres in $\mathbb{R}^n$ are homeomorphic to each other. These examples illustrate that "size" is not a topological property; to a topologist, there is no difference whatsoever between a ball of radius $2$ and a ball of radius $2^{2000}$. As long as one space can be "deformed" into another space by "squeezing and stretching without tearing," the two spaces are identical to the topologists. No wonder these poor souls can't tell the difference between their donut and their coffee mug!*



Figure 2.3: A topologist's perspective.

*b) Let $\mathbb{B}^n \subseteq \mathbb{R}^n$ be the unit ball, and define a map $\Psi : \mathbb{B}^n \to \mathbb{R}^n$ by*

$$\Psi(x) = \frac{x}{1 - |x|}.$$

*A straight computation shows that the map $\Phi \colon \mathbb{R}^n \to \mathbb{B}^n$ defined by*

$$\Phi(x) = \frac{y}{1 + |y|}$$

*is an inverse for $\Psi$. Thus $\Psi$ is bijective, and since $\Psi$ and $\Psi^{-1} = \Phi$ are both continuous, $\Psi$ is a homeomorphism, from which follows that $\mathbb{R}^n$ is homeomorphic to $\mathbb{B}^n$. Thus, since $\mathbb{R}^n$ is unbounded while $\mathbb{B}^n$ is bounded, it follows that "boundedness" is not a topological property either.*

*c) Another illustrative example is the homeomorphism between the surface of a sphere and the surface of a cube: Let $\mathbb{S}^2$ be the unit sphere in $\mathbb{R}^3$, and set*

$$\Omega = \{(x, y, z) \mid \max\{|x|, |y|, |z|\} = 1\},$$

*which is the cubical surface of side 2 centered at the origin. Now let $\varphi \colon \Omega \to \mathbb{S}^2$ be the map that projects each point of $\Omega$ radially inward to the sphere as shown in Figure 2.4.*



Figure 2.4: Deforming the cube $\Omega$ into the unit sphere $\mathbb{S}^2$.

*More precisely, given a point $p \in \Omega$, the image $\varphi(p)$ is the unit vector in the direction of $p$. Thus $\varphi$ is given by the formula*

$$\varphi(x, y, z) = \frac{(x, y, z)}{\sqrt{x^2 + y^2 + z^2}},$$

*which is continuous on $\Omega$ by the usual arguments of elementary analysis (notice that the denominator is always nonzero on $\Omega$, so the mapping is well defined). If we continue to push this argument further, we can show that $\varphi$ is in fact a homeomorphism, with inverse*

$$\varphi^{-1}(x, y, z) = \frac{(x, y, z)}{\max\{|x|, |y|, |z|\}}.$$

*This demonstrates that "corners" are not topological properties either.*

**Definition 69.** *If $X$ is a set, a **basis** for a topology on $X$ is a collection $\mathscr{B}$ of subsets of $X$ (called basis elements) such that*

- *For each $x \in X$, there is at least one basis element $B \in \mathscr{B}$ containing $x$. Note that this says that the basis elements cover the entire space, i.e., $\cup_{B \in \mathscr{B}} = X$.*

- *If $x$ belongs to the intersection of two basis elements $B_1, B_2 \in \mathscr{B}$, then there is another basis element $B_3 \in \mathscr{B}$ containing $x$ such that $B_3 \subset B_1 \cap B_2$.*

*If $\mathscr{B}$ satisfies these two conditions, then we define the **topology $\mathcal{T}$ generated by** $\mathscr{B}$ as follows: A subset $U \subset X$ is said to be open in $X$ (that is, to be an element of $\mathcal{T}$) if for each $x \in U$, there is a basis element $B \in \mathscr{B}$ such that $x \in B$ and $B \subset U$. Note that each basis element is itself an element of $\mathcal{T}$.*

We now define three topologies on the real line $\mathbb{R}$, all of which are of interest to us:

**Definition 70.** *If $\mathscr{B}$ is the collection of all open intervals in the real line,*

$$(a, b) = \{x \mid a < x < b\},$$

*the topology generated by $\mathscr{B}$ is called the **standard topology** on the real line.*

Note that whenever we consider $\mathbb{R}$, unless stated otherwise, we will always assume (as we have always done up until this point) that it is given with this standard topology.

**Definition 71.** *If $\mathscr{B}_\ell$ is the collection of all half-open intervals in the real line,*

$$[a, b) = \{x \mid a \leq x < b\},$$

*where $a < b$, the topology generated by $\mathscr{B}_\ell$ is called the **lower limit topology** on the real line. When $\mathbb{R}$ is given this lower limit topology, we usually denote the space by $\mathbb{R}_\ell$.*

**Example 35.** *Let* $\mathbb{R}$ *denote the set of real numbers in its usual topology, and let* $\mathbb{R}_\ell$ *denote the same set in the lower limit topology. Let* $\mathrm{Id}_\ell \colon \mathbb{R} \to \mathbb{R}_\ell$ *be the identity function* $\mathrm{Id}_\ell(x) = x$ *for every real number* $x$. *Then* $\mathrm{Id}_\ell$ *is not a continuous function; the inverse image of the open set* $[a,b)$ *of* $\mathbb{R}_\ell$ *equals itself, which is not open in* $\mathbb{R}$. *On the other hand, the identity function* $\mathrm{Id}^\ell \colon \mathbb{R}_\ell \to \mathbb{R}$ *is indeed continuous, because the inverse image of* $(a,b)$ *is itself, which is open in* $\mathbb{R}_\ell$. ✧

Example 35 showcases a very important concept that we alluded to earlier in Chapter 1 and I promised we would discuss in this chapter. Notice that $\mathbb{R}_\ell$ is finer than $\mathbb{R}$ (with the usual standard topology), as we will show below in Lemma 18. Since, intuitively, a finer topology contains more open sets than a coarser one, it follows that all maps from a finer topology to a coarser one are continuous (think about this statement and make sure you understand why it's true; it's very important!). This is the reason why we saw in Chapter 1 that all maps from a discrete space are always continuous; it doesn't get any finer than a discrete topology!!

**Definition 72.** *Let* $K = \{1/n \mid n \in \mathbb{N}\}$. *Generate a topology on* $\mathbb{R}$ *by taking as basis* $\mathscr{B}_K$ *the collection of all open intervals* $(a,b)$ *and all the sets of the form* $(a,b) \smallsetminus K$. *The topology generated by* $\mathscr{B}_K$ *is called the* **K-topology** *on the real line. When* $\mathbb{R}$ *is given this topology, we usually denote the space by* $\mathbb{R}_K$.

**Remark:** Relative to the set of all real numbers carrying the *standard topology*, the set $K = \{1/n \mid n \in \mathbb{N}\}$ is not closed since it doesn't contain its (only) limit point 0. Relative to the *K*-topology however, the set $K$ is automatically decreed to be closed by adding more basis elements to the standard topology on $\mathbb{R}$. In other words, the *K*-topology on $\mathbb{R}$ is strictly finer than the standard topology on $\mathbb{R}$, as we show now on Lemma 18.

The relation between these three topologies that we just discussed is the following:

**Lemma 18.** *The topologies of* $\mathbb{R}_\ell$ *and* $\mathbb{R}_K$ *are strictly finer than the standard topology on* $\mathbb{R}$, *but not comparable with one another.*

*Proof.* Let $\mathcal{T}$, $\mathcal{T}_\ell$, and $\mathcal{T}_K$ be the topologies of $\mathbb{R}$, $\mathbb{R}_\ell$, and $\mathbb{R}_K$, respectively. Given any basis element $(a,b)$ for $\mathcal{T}$ and a point $x \in (a,b)$, the basis element $[x,b)$ for $\mathcal{T}_\ell$ has $x$ as an element and moreover, it is contained inside $(a,b)$. On the other hand, given a basis element $[x,d)$ for $\mathcal{T}_\ell$, note that there is no open interval $(a,b)$ that also has $x$ as an element and that is contained inside $[x,d)$. Thus we have shown that $\mathcal{T}_\ell$ is strictly finer than $\mathcal{T}$.

Similarly for $\mathcal{T}_K$, given a basis element $(a, b)$ for $\mathcal{T}$ and a point $x \in (a, b)$, note that this same interval is a basis element for $\mathcal{T}_K$ that also contains $x$. On the other hand, given a basis element $U_K = (-1, 1) \smallsetminus K$ for $\mathcal{T}_K$ and a point $0$ of $U_K$, note that there is no open interval that contains $0$ and lies within $U_K$.

Thus we have shown that the topologies of $\mathbb{R}_\ell$ and $\mathbb{R}_K$ are strictly finer than the standard topology on $\mathbb{R}$. I'll leave it to you as a fun exercise to show that $\mathbb{R}_\ell$ and $\mathbb{R}_K$ are not comparable with one another. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

**Definition 73.** *A **subbasis** $\mathcal{S}$ for a topology on $X$ is a collection of subsets of $X$ whose union equals $X$. The **topology generated by the subbasis** $\mathcal{S}$ is defined to be the collection $\mathcal{T}$ of all unions of finite intersections of elements of $\mathcal{S}$.*

The basic idea is that a basis is the collection of all finite intersections of subbasis elements (i.e. a subbasis generates a basis by taking finite intersections of its elements). Since the open sets in a topology are all possible unions of basis elements, we have that the open sets in a topology are all possible unions of finite intersections of subbasis elements.



Figure 2.5: Basis and subbasis comparison.

**Definition 74.** *Let $X$ be a set with a total order relation (also, assume that $X$ has more than one element). Let $\mathscr{B}$ be the collection of all sets of the following types:*

- *All open intervals $(a, b)$ in $X$.*

- *All intervals of the form $[a_0, b)$, where $a_0$ is the smallest element (if any) of X.*

- *All intervals of the form $(a, b_0]$, where $b_0$ is the largest element (if any) of X.*

*The collection $\mathscr{B}$ is a basis for a topology on X, which is called the **order topology** on X.*

Note that the standard topology on $\mathbb{R}$, as previously defined, is just the order topology derived from the usual order on $\mathbb{R}$.

**Example 36.** *a) The positive integers $\mathbb{Z}^+$ form an ordered set with a smallest element (i.e., 1). The order topology on $\mathbb{Z}^+$ is the discrete topology, for every one-point set is open: If $n > 1$, then the one-point set $\{n\} = (n - 1, n + 1)$ is a basis element; in the case when $n = 1$, the one-point set $\{1\} = [1, 2)$ is a basis element.*

*b): The set $X = \{1, 2\} \times \mathbb{Z}^+$ in the dictionary order is another example of an ordered set with a smallest element. Let $n \in \mathbb{Z}^+$. Then denoting $1 \times n$ by $a_n$ and $2 \times n$ by $b_n$, we can represent X by*

$$a_1, a_2, \ldots; b_1, b_2, \ldots$$

*Note that the order topology on X is not the the discrete topology. The reason is that, even though most one-point sets on X are open, there is an exception –the one-point set $\{b_1\}$. Any open set containing $b_1$ must contain a basis element about $b_1$ (by definition), and any basis element containing $b_1$ contains points of the $a_i$ sequence.* ✹

**Definition 75.** *Let X and Y be topological spaces. The **product topology** on $X \times Y$ is the topology having as basis the collection $\mathscr{B}$ of all sets of the form $U \times V$, where U and V are open subsets of X and Y, respectively.*

**Theorem 73.** *If $\mathscr{B}_X$ is a basis for the topology of X and $\mathscr{B}_Y$ is a basis for the topology of Y, then the collection*

$$\mathscr{B}_{X \times Y} = \{B_X \times B_Y \mid B_X \in \mathscr{B}_X \text{ and } B_Y \in \mathscr{B}_Y\}$$

*is a basis for the topology of $X \times Y$.*

Note that it is sometimes useful to express the product topology in terms of a subbasis. To do this, we first define certain functions called *projections*:

**Definition 76.** *Let $\pi_1 \colon X \times Y \to X$ be defined by the equation*

$$\pi_1(x,y) = x$$

*and let $\pi_2 \colon X \times Y \to Y$ be defined by the equation*

$$\pi_2(x,y) = y.$$

*The maps $\pi_1$ and $\pi_2$ are called the **projections** of $X \times Y$ onto its first and second components, respectively.*

If $U$ is an open subset of $X$, then the set $\pi_1^{-1}(U)$ is precisely the set $U \times Y$, which is open in $X \times Y$. Similarly, if $V$ is open in $Y$, then the set $\pi_2^{-1}(V)$ is precisely the set $X \times V$, which is also open in $X \times Y$. The intersection of these two sets is the set $U \times V$, as indicated in Figure 2.6.



Figure 2.6: $\pi_1^{-1}(U) \cap \pi_2^{-1}(V) = (U \times Y) \cap (X \times V) = U \times V.$

This fact leads us to the following theorem:

**Theorem 74.** *The collection*

$$\mathcal{S} = \left\{ \pi_1^{-1}(U) \mid U \text{ is open in } X \right\} \bigcup \left\{ \pi_2^{-1}(V) \mid V \text{ is open in } Y \right\}$$

*is a subbasis for the product topology on $X \times Y$.*

*Proof.* Let $\mathcal{T}$ denote the product topology on $X \times Y$, and let $\mathcal{T}'$ be the topology generated by $\mathcal{S}$. Because every element of $\mathcal{S}$ belongs to $\mathcal{T}$, so do arbitrary unions of finite intersections of elements of $\mathcal{S}$. Thus $\mathcal{T}' \subset \mathcal{T}$. On the other hand, every basis element $U \times V$ for the topology $\mathcal{T}$ is a finite intersection of elements of $\mathcal{S}$, since

$$U \times V = \pi_1^{-1}(U) \bigcap \pi_2^{-1}(V).$$

Therefore, $U \times V$ belongs to $\mathcal{T}'$, so that $\mathcal{T} \subset \mathcal{T}'$ as well.     $\square$

We now give a more general definition for the product topology, one that we extend to infinite products:

**Definition 77.** *Let $\{X_\alpha\}_{\alpha \in A}$ be an indexed family of topological spaces. Let*

$$\pi_\beta \colon \prod_{\alpha \in A} X_\alpha \to X_\beta$$

*be the projection mapping associated with the index $\beta$, assigning to each element of the product space $\prod_{\alpha \in A} X_\alpha$ its $\beta^{th}$ coordinate,*

$$\pi_\beta((x_\alpha)_{\alpha \in A}) = x_\beta.$$

*Now let $\mathcal{S}_\beta$ denote the collection*

$$\mathcal{S}_\beta = \{\pi_\beta^{-1}(U_\beta) \mid U_\beta \text{ open in } X_\beta\}.$$

*and let $\mathcal{S}$ denote the union of these collections,*

$$\mathcal{S} = \bigcup_{\beta \in A} \mathcal{S}_\beta.$$

*Then the topology generated by the subbasis $\mathcal{S}$ is called the **product topology**. In this topology, $\prod_{\alpha \in A} X_\alpha$ is called the **product space**.*

Now we present another topology that we can define on Cartesian products of topological spaces. As we will see, although the two definitions are strikingly similar, there are alarming differences between the two topologies when dealing with infinite products.

**Definition 78.** *Let $X = \prod_{i \in \mathcal{I}} X_i$, where the $X_i$'s are topological spaces. Then the* **box topology** *on X is the topology generated by the basis*

$$\mathscr{B}_{box} = \left\{ \prod_{i \in \mathcal{I}} U_i \mid U_i \text{ is open in } X_i \right\}.$$

While the definition of the box topology is quite easier to grasp than the more unwieldy definition of the product topology, it is unfortunate that it satisfies fewer desirable properties. In particular, if all the component spaces are compact, the box topology on their Cartesian product will not necessarily be compact, although the product topology on their Cartesian product will always be compact (look up the famous *Tychonoff Theorem* to convince yourself). In general, the box topology is finer than the product topology (all that extra baggage is what messes up all the desirable properties! ☺), although the two agree in the case when our index $\mathcal{I}$ is finite.

**Example 37.** *Take $\mathbb{R}^\omega$, the countable Cartesian product of the real line, and consider the function $f \colon \mathbb{R} \to \mathbb{R}^\omega$ given by*

$$f(x) = (x, x, , x, \ldots);$$

*the $n^{th}$ coordinate function of f is the function $f_n(x) = x$. Each of the coordinate functions $f_n \colon \mathbb{R} \to \mathbb{R}$ is continuous in the standard topology on $\mathbb{R}$, and thus f itself is continuous if $\mathbb{R}^\omega$ is given the product topology, but f is not continuous in the box topology.*

*You wonder why? Well, consider the set*

$$U = \prod_{n=1}^{\infty} (-1/n, 1/n).$$

*This set U is open (it is a basis element) in the box topology, but not in the product topology. We assert that $f^{-1}(U)$ is not open in $\mathbb{R}$. If $f^{-1}(U)$ were open in $\mathbb{R}$, it would contain some interval $(-\delta, \delta)$ about the point 0. But this would mean that $f((-\delta, \delta)) \subset U$, so that, by applying the projection map on the $n^{th}$ coordinate to both sides of this inclusion, we would get*

$$f_n((-\delta, \delta)) = (-\delta, \delta) \subset (-1/n, 1/n) \qquad \text{for all } n.$$

*This is a contradiction because the components of U get arbitrarily close to 0 –any $\delta$-neighborhood will eventually be outside some component of U.*

*The product topology is the more "natural" topology in the sense that pointwise statements about its component spaces are transported to the Cartesian product and vice versa. For instance,*

- *a sequence in the product topology converges if and only if each of the component sequences converges in the component topology,*

- *a function is continuous in the product topology if and only if each of its component functions is continuous in the component topology.*

The following properties of the box topology are very important. Some of these properties may not mean much to you now, but as you proceed with your readings you will run into these concepts later on. Let $\mathbb{R}^\omega$ denote the countable cartesian product of $\mathbb{R}$ with itself. Then the box topology on $\mathbb{R}^\omega$ is:

- completely regular.

- neither compact nor connected.

- not first countable (hence not metrizable).

- not separable.

- paracompact (and hence normal and completely regular) if the continuum hypothesis is true.

**Lemma 19 (The Sequence Lemma).** *Let $X$ be a topological space and let $A \subset X$. If there is a sequence of points of $A$ converging to $x$, then $x \in \bar{A}$. The converse holds if $X$ is metrizable.*

**Example 38.** *Let us show that $\mathbb{R}^\omega$ in the box topology is not metrizable. In order to do this we will show that the Sequence Lemma does not hold for $\mathbb{R}^\omega$. Let $A$ be the subset of $\mathbb{R}^\omega$ consisting of those points all of whose coordinates are positive:*

$$A = \{(x_1, x_2, \dots) \mid x_i > 0 \ \ \forall i \in \mathbb{N}\}.$$

*Let $\tilde{O}$ be the "origin" in $\mathbb{R}^\omega$, that is, the point $(0, 0, \dots)$ each of whose coordinates is zero. In the box topology, $\tilde{O}$ belongs to $\bar{A}$; for if*

$$B = (a_1, b_1) \times (a_2, b_2) \times \dots$$

*is any basis element containing $\widetilde{O}$, then B intersects A. For instance, the point*

$$\left(\frac{1}{2}b_1, \frac{1}{2}b_2, \dots\right)$$

*belongs to $B \cap A$.*

   *But we assert that there is no sequence of points of A converging to $\widetilde{O}$. For let $\{a_n\}$ be a sequence of points of A, where*

$$a_n = (x_{1n}, x_{2n}, \dots, x_{in}, \dots).$$

*Every coordinate $x_{in}$ is positive, so we can construct a basis element $B'$ for the box topology on $\mathbb{R}$ by setting*

$$B' = (-x_{11}, x_{11}) \times (-x_{22}, x_{22}) \times \cdots.$$

*Then $B'$ contains the origin $\widetilde{O}$, but it contains no member of the sequence $\{a_n\}$; the point $a_n$ cannot belong to $B'$ because the $n^{th}$ coordinate $x_{nn}$ does not belong to the interval $(-x_{nn}, x_{nn})$. Hence the sequence $\{a_n\}$ cannot converge to $\widetilde{O}$ in the box topology.*                                      ✇


   In summary, we have the following theorem:


**Theorem 75 (Comparison of the Box and Product Topologies).** *The* box topology *on $\prod X_\alpha$ has as basis all the sets of the form $\prod U_\alpha$, where $U_\alpha$ is open in $X_\alpha$ for each $\alpha$. The* product topology *on $\prod X_\alpha$ has as basis all the sets of the form $\prod U_\alpha$, where $U_\alpha$ is open in $X_\alpha$, **and** $U_\alpha$ equals $X_\alpha$ except for finitely many values of $\alpha$.*


**Theorem 76.** *Let $\{X_\alpha\}$ be an indexed family of spaces, and let $A_\alpha \subset X_\alpha$ for each $\alpha$. If $\prod X_\alpha$ is given by either the product or the box topology, then we have*

$$\prod \overline{A_\alpha} = \overline{\prod A_\alpha}.$$


**Definition 79.** *Let X be a topological space with topology $\mathcal{T}$. If Y is a subset of X, then the collection*

$$\mathcal{T}_Y = \{Y \cap U \mid U \in \mathcal{T}\}$$

*is a topology on Y, called the **subspace topology**. Equipped with this topology, Y is called a **subspace** of X.*

**Lemma 20.** *If $\mathscr{B}_X$ is a basis for the topology of X, then the collection*

$$\mathscr{B}_Y = \{B \cap Y \mid B \in \mathscr{B}_X\}$$

*is a basis for the subspace topology on Y.*

Now let us explore the relation between the subspace, order, and product topologies. For product topologies, the result is what one might expect; for order topologies however, things are not so hunky dory:

**Theorem 77.** *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 $X \times Y$.*

*Proof.* The set $U \times V$ is the general basis element for $X \times Y$, where $U$ is open in $X$ and $V$ is open in $Y$. Therefore $(U \times V) \cap (A \times B)$ is the general basis element for the subspace topology on $A \times B$. Now

$$(U \times V) \cap (A \times B) = (U \cap A) \times (V \cap B).$$

Since $U \cap A$ and $V \cap B$ are general open sets for the subspace topologies on $A$ and $B$, respectively, the set $(U \cap A) \times (V \cap B)$ is the general basis element for the product topology on $A \times B$.

The conclusion we draw is that the bases for the subspace topology on $A \times B$ and the product topology on $A \times B$ are the same. Hence the topologies are the same. □

Now let $X$ be an ordered set in the order topology, and let $Y$ be a subset of $X$. The order relation on $X$, when restricted to $Y$, makes $Y$ into an ordered set. However, the resulting order topology on $Y$ need not be the same as the topology that $Y$ inherits as a subspace of $X$.

**Example 39.** *a) We first show an example where the subspace and order topologies on Y agree: Consider the subset $Y = [0,1] \subset \mathbb{R}$, in the subspace topology. The subspace topology has as basis all sets of the form $(a,b) \cap Y$, where $(a,b)$ is an open interval in $\mathbb{R}$. Such a set is one of the following types:*

$$(a,b) \cap Y = \begin{cases} (a,b) & \text{if } a \text{ and } b \text{ are in } Y, \\ [0,b) & \text{if only } b \text{ is in } Y, \\ (a,1] & \text{if only } a \text{ is in } Y, \\ Y \text{ or } \varnothing & \text{if neither } a \text{ nor } b \text{ is in } Y. \end{cases}$$

*By definition, each of these sets is open in Y, but the ones of type $[0, b)$ and $(a, 1]$ are not open in the ambient space $\mathbb{R}$. Note that these sets form a basis for the order topology on Y. Thus, we see that in the case of the set $Y = [0, 1]$, its subspace topology (as a subspace of $\mathbb{R}$) and its order topology are the same.*

*b) We now show an example where the subspace and order topologies on Y don't agree: Let Y be the subset $[0, 1) \cup \{2\}$ of $\mathbb{R}$. In the subspace topology on Y, the one-point set $\{2\}$ is open, because it is the intersection of the open set $(3/2, 5/2)$ with Y. But in the order topology on Y, the set $\{2\}$ is not open. Any basis element for the order topology on Y that contains 2 is of the form*

$$\{x \mid x \in Y \text{ and } a < x \leq 2\}$$

*for some $a \in Y$. Such a set necessarily contains points of Y less than 2. Hence we see that the subspace and order topologies on Y do not agree.*

We have however, that if $Y$ is convex in $X$, then the subspace and order topologies do always agree:

**Theorem 78.** *Let X be an ordered set in the order topology, and let Y be a subset of X that is convex in X. Then the order topology on Y is the same as the topology Y inherits as a subspace of X.*

*Proof.* $(\subseteq)$ Consider the ray $(a, \infty)$ in $X$. If $a \in Y$, then we have

$$(a, \infty) \bigcap Y = \{x \mid x \in Y, \, x > a\},$$

which is an open ray of the ordered set $Y$. Now, if $a \notin Y$, then $a$ is either a lower bound or an upper bound on $Y$. If it's a lower bound, then note that $(a, \infty) \cap Y = Y$; if instead $a$ is an upper bound, then $(a, \infty) \cap Y = \emptyset$.

A similar argument shows that $(-\infty, a) \cap Y$ is either an open ray of $Y$, all of $Y$ itself, or empty. Since the sets $(a, \infty) \cap Y$ and $(-\infty, a) \cap Y$ form a subbasis for the subspace topology on $Y$, and since each of them is open in the order topology, it follows that the order topology contains the subspace topology.

$(\supseteq)$ To show the reverse containment, note that any open ray of $Y$ equals the intersection of an open ray of $X$ with $Y$, so it is open in the subspace topology on $Y$. Since the open rays of $Y$ are a subbasis for the order topology on $Y$, this topology is contained in the subspace topology, as desired. $\square$

**Theorem 79.** *Let $Y$ be a subspace of $X$. Then a set $A$ is closed in $Y$ if and only if it equals the intersection of a closed set of $X$ with $Y$.*

**Example 40.** *Consider the subspace $Y = (0,1]$ of the real line $\mathbb{R}$. The set $A = (0,1/2)$ is a subset of $Y$; its closure in $\mathbb{R}$ is the set $[0,1/2]$, while its closure in $Y$ is the set $[0,1/2] \cap Y = (0,1/2]$.* ✇

### 2.1.1   Hausdorff Condition

**Definition 80.** *A topological space $X$ is said to be a **Hausdorff space** (also a $T_2$ **space**)[1] if given any pair of distinct points $p_1, p_2 \in X$, there exist neighborhoods $U_1$ of $p_1$ and $U_2$ of $p_2$ with $U_1 \cap U_2 = \emptyset$. (This property is often summarized by saying that "points can be separated by open subsets.")*

**Example 41.** *a) Every metric space is Hausdorff: if $p_1$ and $p_2$ are distinct, let $r = d(p_1, p_2)$; then the open balls of radius $r/2$ around $p_1$ and $p_2$ are disjoint by the triangle inequality.*

*b) Every discrete space is Hausdorff, because $\{p_1\}$ and $\{p_2\}$ are disjoint open subsets when $p_1 \neq p_2$.*

*c) Every open subset of a Hausdorff space is Hausdorff: if $V \subseteq X$ is open in the Hausdorff space $X$, and $p_1, p_2$ are distinct points in $V$, then in $X$ there are open subsets $U_1, U_2$ separating $p_1$ and $p_2$, and the sets $U_1 \cap V$ and $U_2 \cap V$ are open in $V$, disjoint, and contain $p_1$ and $p_2$, respectively.*

*d) Suppose $X$ is a topological space, and for every $p \in X$ there exists a continuous function $f : X \to \mathbb{R}$ such that $f^{-1}(\{0\}) = \{p\}$. It can be shown that $X$ is Hausdorff.*

*e) The trivial topology on any set containing more than one element is **not** Hausdorff, nor is the topology on $\{1,2,3\}$ described in* Example 33 c). *Because every metric space is Hausdorff, it follows that these spaces are not metrizable.*

Hausdorff spaces have many of the properties that we expect of metric spaces, such as those expressed in the following proposition:

**Proposition 26.** *Let $X$ be a Hausdorff space.*

✦────────────── ☞

---

[1] This terminology comes from the *separation axioms*, which we will discuss in greater detail later on. A $T_1$ *space* will also be defined shortly.

a) *Every finite subset of X is closed.*

b) *If a sequence $\{(p_i\}$ in X converges to a limit $p \in X$, the limit is unique.*

*Proof of a).* Consider first a set $\{p_0\}$ containing only one point. Given $p \neq p_0$, the Hausdorff property says that there exist disjoint neighborhoods $U$ of $p$ and $V$ of $p_0$. In particular, $U$ is a neighborhood of $p$ contained in $X \smallsetminus \{p_0\}$, so $\{p_0\}$ is closed. It then follows that finite subsets are closed, because they are simply finite unions of one-point sets and, as we already know, finite unions of closed sets are closed. $\qquad\square$

*Proof of b).* Suppose, to the contrary, that a sequence $\{p_i\}$ has two distinct limits $p$ and $p'$. By the Hausdorff property, there exist disjoint neighborhoods $U$ of $p$ and $U'$ of $p'$. Also, by the definition of convergence, there exist $N, N' \in \mathbb{N}$ such that $i \geq N$ implies $p_i \in U$ and $i \geq N'$ implies $p_i \in U'$. But since $U$ and $U'$ are disjoint, this is a contradiction when $i \geq \max\{N, N'\}$. $\qquad\square$

Note, however, that the condition that a finite point set is closed does not guarantee 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. This condition that finite point sets be closed has a name of its own: it is called the $T_1$ *axiom*, which we now define:

**Definition 81.** *We say that a space X is a $T_1$ **space** if for any two points $x, y \in X$ there exists two open sets U and V such that $x \in U$ and $y \notin U$, and $y \in V$ and $x \notin V$.*

While we're at it, we might as well define another one of these weird separation axioms:

**Definition 82.** *We say that a space X is a $T_0$ **space** (or a **Kolmogorov space**) if for any two points $x, y \in X$, there is an open set U such that $x \in U$ and $y \notin U$ or $y \in U$ and $x \notin U$.*

Alright, that's enough for now; we will see more of these weirdos on Section 2.2. Moving on, it can be shown that the only Hausdorff topology on a finite set is the discrete topology. Another important property of Hausdorff spaces is expressed in the following proposition:

**Proposition 27.** *Suppose X is a Hausdorff space and $A \subseteq X$. If $p \in X$ is a limit point of A, then every neighborhood of p contains infinitely many points of A.*

And here's the analogous result for $T_1$ spaces, which is a biconditional statement:

**Proposition 28.** *Let X be a space satisfying the $T_1$ axiom and let A be a subset of X. Then the point x is a limit point of A if and only if every neighborhood of x contains infinitely many points of A.*

We close out this subsection by stating two important results related to Hausdorff spaces.

**Theorem 80.** *Every totally 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.*

**Definition 83.** *A map $f : X \to Y$ is said to be a **closed map** if for each closed set A of X, the set $f(A)$ is closed in Y. Similarly, $f : X \to Y$ is said to be an **open map** if for each open set U of X, the set $f(U)$ is open in Y.*

**Lemma 21 (Closed Map Lemma).** *Suppose F is a continuous map from a compact space to a Hausdorff space. Then,*

    *a)* *F is a closed map.*

    *b)* *If F is surjective, it is a quotient map.*

    *c)* *If F is injective, it is a topological embedding.*

    *d)* *If F is bijective, it is a homeomorphism.*

### 2.1.2   Quotient Topology

**Definition 84.** *Let X and Y be topological spaces, and let $p : X \to Y$ be a surjective map. The map p is said to be a **quotient map** provided a subset U of Y is open in Y if and only if $p^{-1}(U)$ is open in X.*

Note that, by its biconditional nature, this concept of a quotient map is stronger than that of a continuous map; some mathematicians call this condition "strong continuity." An equivalent formulation is to require that a subset $A$ of $Y$ be closed if and only if $p^{-1}(A)$ is closed in $X$. Equivalence of the two statements follow from the fact that $f^{-1}(Y \smallsetminus B) = X \smallsetminus f^{-1}(B)$.

Yet another way of describing a quotient map is as follows: We say that a subset $C$ of $X$ is **saturated** (with respect to the surjective map $p \colon X \to Y$) if $C$ contains every set $p^{-1}(\{y\})$ that it intersects. In other words, $C$ is saturated if it equals the complete inverse image of a subset of $Y$. Thus, to say that $p$ is a quotient map is equivalent to saying that $p$ is continuous and $p$ maps saturated open sets of $X$ to open sets of $Y$ (or saturated closed sets of $X$ to closed sets of $Y$).

Two special kinds of quotient maps are the open maps and the closed maps. It follows immediately from the definition that if $p \colon X \to Y$ is a surjective continuous map that is either open or closed, then $p$ is a quotient map. The converse however does not always hold; that is, there are quotient maps that are neither open nor closed.

**Example 42.** *a) Let $X$ be the subspace $[0,1] \cup [2,3]$ of $\mathbb{R}$ and let $Y$ be the subspace $[0,2]$ of $\mathbb{R}$. The map $p \colon X \to Y$ defined by*

$$p(x) = \begin{cases} x & \text{for } x \in [0,1], \\ x - 1 & \text{for } x \in [2,3] \end{cases}$$

*is readily seen to be surjective, continuous, and closed. Therefore it is a quotient map. It is not, however, an open map; the image of the open set $[0,1]$ of $X$ is not open in $Y$. Note that if $A$ is the subspace $[0,1) \cup [2,3]$ of $X$, then the map $q \colon A \to Y$ obtained by restricting $p$ is continuous and surjective, but it is not a quotient map. The reason is that the set $[2,3]$ is open in $A$ and is saturated with respect to $q$, but its image is not open in $Y$.*

*b) Let $\pi_1 \colon \mathbb{R} \times \mathbb{R} \to \mathbb{R}$ be the projection onto the first coordinate; then $\pi_1$ is continuous and surjective. Furthermore, $\pi_1$ is an open map. For if $U \times V$ is a nonempty basis element for $\mathbb{R} \times \mathbb{R}$, then $\pi_1(U \times V) = U$ is open in $\mathbb{R}$; it follows that $\pi_1$ carries open sets of $\mathbb{R} \times \mathbb{R}$ to open sets of $\mathbb{R}$. However, $\pi_1$ is not a closed map. The subset $C = \{x \times y \mid xy = 1\} \subset \mathbb{R} \times \mathbb{R}$ is closed, but $\pi_1(C) = \mathbb{R} \smallsetminus \{0\}$, which is not closed in $\mathbb{R}$. Note that if $A \subset \mathbb{R} \times \mathbb{R}$ is such that $A = C \cup \{0\}$, then the map $q \colon A \to \mathbb{R}$ obtained by restricting $\pi_1$ is continuous and surjective, but it is not a quotient map. The reason is that the one-point set $\{0\}$ is open in $A$ and is saturated with respect to $q$, but its image is not open in $\mathbb{R}$.*

We now show how the notion of a quotient map can be used to construct a topology on a set:

**Definition 85.** *If X is a space and S is a set and if $p\colon X \to S$ is a surjective map, then there exists exactly one topology $\mathcal{T}$ on S relative to which p is a quotient map; it is called the **quotient topology induced by** p.*

The topology $\mathcal{T}$ mentioned in the above definition is of course defined by letting it consist of those subsets $U$ of $S$ such that $p^{-1}(U)$ is open in $X$. It is easy to check that $\mathcal{T}$ is a topology:

- The sets $\varnothing$ and $S$ are open because

$$p^{-1}(\varnothing) = \varnothing \qquad \text{and} \qquad p^{-1}(S) = X.$$

- The other two conditions follow from the equations

$$p^{-1}\left(\bigcup_{\alpha \in A} U_\alpha\right) = \bigcup_{\alpha \in A} p^{-1}(U_\alpha)$$

and

$$p^{-1}\left(\bigcap_{i=1}^{n} U_i\right) = \bigcap_{i=1}^{n} p^{-1}(U_i).$$

**Example 43.** *Let p be the map of the real line $\mathbb{R}$ onto the three-point set $S = \{a, b, c\}$ defined by*

$$p(x) = \begin{cases} a & \text{if } x > 0, \\ b & \text{if } x < 0, \\ c & \text{if } x = 0. \end{cases}$$

*You can check that the quotient topology on S induced by p is the one indicated in Figure 2.7.*  ✦

There is a special situation in which the quotient topology occurs particularly frequently. It is the following:

Figure 2.7: Topology on $S$ induced by $p$.

**Definition 86.** *Let $X$ be a topological space, and let $X^*$ be a partition of $X$ into disjoint subsets whose union is $X$. Let $p\colon X \to X^*$ be the surjective map that carries each point of $X$ to the element of $X^*$ containing it. In the quotient topology induced by $p$, the space $X^*$ is called the* **quotient space** *of $X$.*

**Example 44.** *Let $X$ be the closed unit ball*

$$\mathbb{B}^2 = \{(x,y) \mid x^2 + y^2 \leq 1\}$$

*in $\mathbb{R}^2$, and let $X^*$ be the partition of $X$ consisting of all the one-point sets $\{(x,y)\}$ for which $x^2 + y^2 < 1$, along with the circle $\mathbb{S}^1 = \{(x,y) \mid x^2 + y^2 = 1\}$. Typical saturated open sets in $X$ are pictured by the shaded regions in Figure 2.8. One can show that $X^*$ is homeomorphic to the 2-sphere $\mathbb{S}^2 \subset \mathbb{R}^3$, given by $\mathbb{S}^2 = \{(x,y,z) \mid x^2 + y^2 + z^2 = 1\}$.*



Figure 2.8: Topology on $A$ induced by $p$.

## 2.2   Countability & Separation Axioms

In this section we will present the famous separation axioms of general topology after briefly presenting the axioms of countability, of which the *second axiom* is of most importance for most applications. For instance, as you will see in Chapter 4, the second countability axiom is a requirement that we impose on our definition of manifolds. Some authors don't impose this requirement from the get-go, but then they end up using unnecessarily long phrases such as "Let $M$ be a Hausdorff, second countable manifold..." By requiring that manifolds be Hausdorff and second countable from the start, we are guaranteed the existence of a compact exhaustion and a partition of unity for our manifold, which are highly desirable tools to possess for most applications. These words may not mean much to you at the moment, but as always I present them to you upfront to spark your curiosity and so that when you get to see the heavy machinery later on, you will get to appreciate it more.

### 2.2.1   Countability Axioms

**Definition 87.** *A space $X$ is said to have a **countable basis at a point** $x \in X$ if there is a countable collection $\mathscr{B}$ of neighborhoods of $x$ such that each neighborhood of $x$ contains at least one of the elements of $\mathscr{B}$. A space that has a countable basis at each of its points is said to satisfy the **first countability axiom**, or to be **first-countable**.*

Note that every metrizable space satisfies this axiom.

**Theorem 81.** *Let $X$ be a topological space. Then,*

- *let $A$ be a subset of $X$. If there is a sequence of points of $A$ converging to $x$, then $x \in \bar{A}$. The converse holds if $X$ is first-countable.*

- *let $f \colon X \to Y$. If $f$ is continuous, for every convergent sequence $x_n \to x$ in $X$, the sequence $f(x_n)$ converges to $f(x)$. The converse holds if $X$ is first-countable.*

As I emphasized at the beginning of this section, of much greater importance than the first countability axiom is the following:

**Definition 88.** *If a space X has a countable basis for its topology, then X is said to satisfy the* **second countability axiom**, *or to be* **second-countable**.

It is quite obvious that the second axiom implies the first: Let $\mathscr{B} = \{U_n\}$ be a countable basis for the topology of a space $X$. For each point $x \in X$, the subcollection $\mathscr{B}_x = \{U \in \mathscr{B} \mid x \in U\}$ is a countable neighborhood basis at $x$.

**Proposition 29.** *We have the following two crucial properties of first and second countable spaces:*

   a) *A subspace of a first countable (resp. second countable) space is first countable (resp. second countable).*

   b) *A countable product of first countable (resp. second countable) spaces is first countable (resp. second-countable).*

We will prove the second countable cases; the first countable counterparts follow a similar argument.

*Proof of a).* If $\mathscr{B} = \{U_n\}$ is a countable basis for $X$ and $A \subset X$, then $\{A \cap U_n\}$ is a countable basis for $A$. □

*Proof of b).* If $\mathscr{B}_n$ is a countable basis for $X_n$, where $n \in \mathbb{Z}^+$, then the collection of products

$$\prod_{n=1}^{\infty} U_n,$$

where $U_n \in \mathscr{B}_n$ for finitely many $n$ and $U_n = X_n$ for the remaining $n$ is a countable basis for the product topology on $\prod_{n=1}^{\infty} X_n$. □

**Theorem 82.** *Suppose that X has a countable basis.*

   • *Every open covering of X contains a countable subcovering X. A space satisfying this condition is called a* **Lindelöf space**.

   • *There exists a countable subset of X that is dense in X. In this case the space X is said to be* **separable**.

The two properties stated in the previous theorem are generally weaker than the second countability axiom, albeit they are equivalent when the space $X$ is metrizable. Other relevant countability axioms for topological spaces include:

- *sequential space:* a set $X$ is open if every sequence convergent to a point in the set is eventually in $X$.

- *$\sigma$-compact space:* there exists a countable cover by compact spaces.

Lastly, here's a summary of the relationships amongst all the countability axioms we have discussed:

- Every first countable space is sequential.

- Every $\sigma$-compact space is Lindelöf.

- Every second countable space is first countable, separable, and Lindelöf.

- Every metric space is first countable.

- For metric spaces, second countability, separability, and the Lindelöf property are all equivalent.

### 2.2.2   Separation Axioms

Throughout our discussions of the separation axioms, we will assume that all spaces being considered satisfy the $T_1$ axiom, i.e., that all finite point sets (in particular one-point sets) are closed. Otherwise we would run into some pathological cases that will be more trouble than they're worth. To see an example of why we want this axiom satisfied, see the remark immediately after the following two definitions.

**Definition 89.** *Suppose that one-point sets are closed in $X$. Then $X$ is said to be* **regular** *if for each pair consisting of a point $x$ and a closed set $B$ disjoint from $x$, there exists disjoint open sets containing $x$ and $B$, respectively. A space that is regular (and satisfies the $T_1$ axiom) is said to satisfy the $T_3$* **axiom***, or to be a $T_3$* **space***.*

**Definition 90.** *Suppose that one-point sets are closed in X. Then X is said to be **normal** if for each pair A, B of disjoint closed sets of X, there exist disjoint open sets containing A and B, respectively. A space that is normal (and satisfies the $T_1$ axiom) is said to satisfy the $T_4$ **axiom**, or to be a $T_4$* ***space.***

**Remark:** It is clear that a regular space is Hausdorff, and that a normal space is regular. However, note that we need to include the condition that one-point sets be closed (i.e., the $T_1$ axiom) as part of the definition of regularity and normality in order for this to be the case. To see why, take for instance a two-point space $Y$ in the indiscrete topology. Then $Y$ satisfies the other part of the definitions of regularity and normality, even though it is not Hausdorff. To see examples showing that regularity is stronger than Hausdorff, and normality is stronger than regularity, check Example 45 below.

The three separation axioms are illustrated in Figure 2.9.



Figure 2.9: Normal $\implies$ Regular $\implies$ Hausdorff.

There is actually a few more separation axioms besides these three and the previously mentioned $T_0$, but we will not discuss those in this section; we're moving on to the fun stuff (aka. algebraic topology) right after this section! Anyhow, there are other ways to formulate the separation axioms that we have presented here. One such formulation that is sometimes useful is given by the following lemma:

**Lemma 22.** *Let X be a topological space and let one-point sets in X be closed.*

a) *X is regular if and only if given a point $x \in X$ and a neighborhood $U$ of $x$, there is a neighborhood $V$ of $x$ such that $\overline{V} \subset U$.*

b) *X is normal if and only if given a closed set $A$ and an open set $U$ containing $A$, there is an open set $V$ containing $A$ such that $\overline{V} \subset U$.*

*Proof of a).* ($\Rightarrow$) Suppose that $X$ is regular, and suppose that the point $x$ and the neighborhood $U$ of $x$ are given. Let $B = X \smallsetminus U$, so that $B$ is a closed set. By hypothesis, there exists disjoint open sets $V$ and $W$ containing $x$ and $B$, respectively. The set $\overline{V}$ is disjoint from $B$, since if $y \in B$, the set $W$ is a neighborhood of $y$ disjoint from $V$. Therefore, $\overline{V} \subset U$, as desired.

($\Leftarrow$) To prove the converse, suppose the point $x$ and the closed set $B$ not containing $x$ are given. Let $U = X \smallsetminus B$. By hypothesis, there is a neighborhood $V$ of $x$ such that $\overline{V} \subset U$. The open sets $V$ and $X \smallsetminus \overline{V}$ are disjoint open sets containing $x$ and $B$, respectively. Thus, $X$ is regular, as desired. $\qquad\square$

*Proof of b).* This proof uses exactly the same argument as in part *a*). We only need to replace the point $x$ by the set $A$ throughout. $\qquad\square$

Now we relate the separation axioms with the concepts previously introduced:

**Theorem 83.** *We have the following properties:*

a) *A subspace of a Hausdorff space is Hausdorff. Also, a product of Hausdorff spaces is Hausdorff.*

b) *A subspace of a regular space is regular. Also, a product of regular spaces is regular.*

*Proof of a).* Let $X$ be Hausdorff. Let $x$ and $y$ be two points of the subspace $Y$ of $X$. If $U$ and $V$ are disjoint neighborhoods in $X$ of $x$ and $y$, respectively, then $U \cap Y$ and $V \cap Y$ are disjoint neighborhoods of $x$ and $y$ in $Y$. Hence, a subspace of a Hausdorff space is Hausdorff.

Now let $\{X_\alpha\}$ be a family of Hausdorff spaces. Let $x = (x_\alpha)$ and $y = (y_\alpha)$ be distinct points of the product space $\prod X_\alpha$. Because we are assuming that $x \neq y$, there is some index $\beta$ such that $x_\beta \neq y_\beta$. Choose disjoint open sets $U$ and $V$ in $X_\beta$ containing $x_\beta$ and $y_\beta$, respectively. Then the sets $\pi_\beta^{-1}(U)$ and $\pi_\beta^{-1}(V)$ are disjoint open sets in $\prod X_\alpha$ containing $x$ and $y$, respectively. Thus, a product of Hausdorff spaces is Hausdorff, as desired. $\qquad\square$

*Proof of b).* Let $Y$ be a subspace of the regular space $X$. Then one-point sets are closed in $Y$. Let $x$ be a point of $Y$ and let $B$ be a closed subset of $Y$ disjoint from $x$. Now $\overline{B} \cap Y = B$, where $\overline{B}$ denotes the closure of $B$ in $X$. Therefore $x \notin \overline{B}$, so using regularity of $X$, we can choose disjoint open sets $U$ and $V$ of $X$ containing $x$ and $\overline{B}$, respectively. Then $U \cap Y$ and $V \cap Y$ are disjoint open sets in $Y$ containing $x$ and $B$, respectively. Hence, we have that a subspace of a regular space is regular.

Now let $\{X_\alpha\}$ be a family of regular spaces, and let $X = \prod X_\alpha$. By part $a$), $X$ is Hausdorff, so that one-point sets are closed in $X$. We now use Lemma 22 to prove regularity of $X$: Let $x = (x_\alpha)$ be a point of $X$ and let $U$ be a neighborhood of $x$ in $X$. Choose a basis element $\prod U_\alpha$ about $x$ contained in $U$. Choose, for each $\alpha$, a neighborhood $V_\alpha$ of $x_\alpha$ in $X_\alpha$ such that $\overline{V}_\alpha \subset U_\alpha$; if it happens that $U_\alpha = X_\alpha$, then choose $V_\alpha = X_\alpha$. Then $V = \prod V_\alpha$ is a neighborhood of $x$ in $X$. Since $\overline{V} = \prod \overline{V}_\alpha$ by Theorem 76, it follows at once that $\overline{V} \subset \prod U_\alpha \subset U$, so that $X$ is regular. Thus, we have shown that a product of regular spaces is regular, and we are done. $\qquad\qquad\square$

There is no analogous theorem for normal spaces, as we shall see shortly. First, a theorem that we are going to use on the next example:

**Theorem 84.** *Let $A$ be a set. There is no injective map $f\colon \mathcal{P}(A) \to A$ (where $\mathcal{P}(-)$ denotes as usual the power set), and there is no surjective map $g\colon A \to \mathcal{P}(A)$.*

**Example 45.** *a) The space $\mathbb{R}_K$ is Hausdorff but not regular. It is Hausdorff because any two distinct points have disjoint open intervals containing them. However it is not regular, for the set $K$ is closed in $\mathbb{R}_K$ and it does not contain the point $\{0\}$: Suppose that there exist disjoint open sets $U$ and $V$ containing $\{0\}$ and $K$, respectively. Now choose a basis element containing $\{0\}$ and lying in $U$. It must be a basis element of the form $(a,b) \smallsetminus K$, since each basis element of the form $(a,b)$ containing $\{0\}$ intersects $K$. Choose $n$ large enough so that $1/n \in (a,b)$. Then choose a basis element about $1/n$ contained in $V$; it must be a basis element of the form $(c,d)$. Finally, choose $z$ such that $\max\{c, 1/(n+1)\} < z < 1/n$. Then $z$ belongs to both $U$ and $V$, so they are not disjoint.*

*b) The space $\mathbb{R}_\ell$ is normal. It is immediate that one-point sets are closed in $\mathbb{R}_\ell$, since the topology of $\mathbb{R}_\ell$ is finer than that of $\mathbb{R}$. To check normality, suppose that $A$ and $B$ are disjoint closed sets in $\mathbb{R}_\ell$. For each point $a \in A$, choose a basis element $[a, x_a)$ not intersecting $B$; and for each point $b \in B$, choose a basis element $[b, x_b)$ not intersecting $A$. The open sets*

$$U = \bigcup_{a \in A} [a, x_a) \qquad and \qquad V = \bigcup_{b \in B} [b, x_b)$$

*are disjoint open sets of A and B, respectively.*

*c) The space $\mathbb{R}_\ell^2$ (i.e., $\mathbb{R}^2$ endowed with the lower limit topology), which is commonly referred to as the **Sorgenfrey plane**, is not a normal space. Note that the space $\mathbb{R}_\ell$ is regular (in fact, normal, as we saw in part b)), so the product space $\mathbb{R}_\ell^2$ is also regular. Hence this example serves two purposes: it shows that a regular space need not be normal, and it shows that the product of two normal spaces need not be normal: We suppose that $\mathbb{R}_\ell^2$ is normal and then arrive at a contradiction. Let $\Delta$ be the diagonal, i.e., the subspace of $\mathbb{R}_\ell^2$ consisting of all points of the form $x \times (-x)$.*



$\Delta = \{\, (x, -x) \mid x \in \mathbb{R} \,\}$

*Then $\Delta$ is closed in $\mathbb{R}_\ell^2$ and it has the discrete topology. Hence every subset A of $\Delta$, being closed in $\Delta$, is also closed in $\mathbb{R}_\ell^2$. Because $\Delta \smallsetminus A$ is also closed in $\mathbb{R}_\ell^2$, this means that for every nonempty proper subset A of $\Delta$, one can find disjoint open sets $U_A$ and $V_A$ containing A and $\Delta \smallsetminus A$, respectively.*

*Let D denote the set of points of $\mathbb{R}_\ell^2$ having rational coordinates; it is dense in $\mathbb{R}_\ell^2$. We define a map $\Psi$ that assigns, to each subset of the the line $\Delta$, a subset of the set D, by setting*

$$\Psi(A) = D \bigcap U_A \qquad \text{if } \varnothing \subsetneq A \subsetneq \Delta.$$
$$\Psi(\varnothing) = \varnothing,$$
$$\Psi(\Delta) = D.$$

*We now show that $\Psi : \mathcal{P}(\Delta) \to \mathcal{P}(D)$ is injective: Let A be a proper nonempty subset of $\Delta$. Then $\Psi(A) = D \cap U_A$ is neither empty (since $U_A$ is open and D is dense in $\mathbb{R}_\ell^2$) nor all of D (since $D \cap V_A$ is nonempty).*

*It remains to show that if B is another proper nonempty subset of $\Delta$, then $\Psi(A) \neq \Psi(B)$. One of the sets A, B contains a point not in the other; suppose WLOG that $x \in A$ and $x \notin B$. Then $x \in \Delta \smallsetminus B$, so that $x \in U_A \cap V_B$; since the latter set is open and nonempty, it must contain points of D. These points belong to $U_A$ and not to $U_B$; therefore, $D \cap U_A \neq D \cap U_B$, as desired. Thus, $\Psi$ is injective.*

*Now we show that there exists an injective map $\mathcal{P}(D) \to \Delta$. Because $D$ is countably infinite and $\Delta$ has the cardinality of $\mathbb{R}$, it suffices to define an injective map $\psi$ of $\mathcal{P}(\mathbb{Z}^+)$ into $\mathbb{R}$. For that, we let $\psi$ assign to the subset $S$ of $\mathbb{Z}^+$ the infinite decimal $.a_1 a_2 \ldots$, where*

$$a_i = \begin{cases} 0 & \text{if } i \in S, \\ 1 & \text{if } i \notin S. \end{cases}$$

*That is,*

$$\psi(S) = \sum_{i=1}^{\infty} \frac{a_i}{10^i}.$$

*Now the composite*

$$\mathcal{P}(\Delta) \xrightarrow{\Psi} \mathcal{P}(D) \xrightarrow{\psi} \Delta$$

*is an injective map of $\mathcal{P}(\Delta)$ into $\Delta$. But Theorem 84 tells us that such a map does not exist! Thus we have reached the desired contradiction.*

**Theorem 85.** *Every regular space with a countable basis is normal.*

*Proof.* Let $X$ be a regular space with a countable basis $\mathcal{B}$. Let $A$ and $B$ be disjoint closed subsets of $X$. Each point $x \in A$ has a neighborhood $U$ not intersecting $B$. Using regularity, choose a neighborhood $V$ of $x$ whose closure lies in $U$, and then choose an element of $\mathcal{B}$ containing $x$ and contained in $V$. By choosing such a basis element for each $x \in A$, we construct a countable covering of $A$ by open sets whose closures do not intersect $B$. Since this covering of $A$ is countable, we can index it with the positive integers; let us denote it by $\{U_n\}$.

Similarly, choose a countable collection $\{V_n\}$ of open sets covering $B$, such that each set $\overline{V}_n$ is disjoint from $A$. The sets $U = \cup U_n$ and $V = \cup V_n$ are open sets containing $A$ and $B$, respectively, but they need not be disjoint. Thus, in order to construct two open sets that are disjoint, we perform the following simple trick: Given $n$, define

$$U_n' = U_n \setminus \bigcup_{i=1}^{n} \overline{V}_i \quad \text{and} \quad V_n' = V_n \setminus \bigcup_{i=1}^{n} \overline{U}_i.$$

Note that each set $U_n'$ is open, being the difference of an open set $U_n$ and a closed set $\cup_{i=1}^{n} \overline{V}_i$. Similarly, each set $V_n'$ is open. It follows that the collection $\{U_n'\}$ covers $A$, because

Figure 2.10: Proof of Theorem 85.

each $x \in A$ belongs to $U_n$ for some $n$, and $x$ belongs to none of the sets $\overline{V}_i$. Similarly, the collection $\{V'_n\}$ covers $B$. (See Figure 2.10.)

Finally, the open sets

$$U' = \bigcup_{n \in \mathbb{Z}^+} U'_n \quad \text{and} \quad V' = \bigcup_{n \in \mathbb{Z}^+} V'_n$$

are disjoint. For if $x \in U' \cap V'$, then $x \in U'_j \cap V'_k$ for some $j$ and $k$. Suppose that $j \leq k$. It then follows from the definition of $U'_j$ that $x \in U_j$; and since $j \leq k$, it follows from the definition of $V'_k$ that $x \notin \overline{U}_j$. A similar contradiction arises if $j \geq k$. $\qquad\square$

**Theorem 86.** *Every metrizable space is normal.*

**Theorem 87.** *Every compact Hausdorff space is normal.*

**Theorem 88.** *Every well-ordered set X is normal in the order topology.*

To see proofs of these three theorems, and to see related examples, see [Munkres, 2000, p. 202-205].

## 2.3 Basics of Algebraic Topology

### 2.3.1 Some Categorical Nonsense

Mathematics is often referred to as the language of the universe. Well, as it turns out, it could be argued that category theory is (in a sense) the universal language of mathematics. Just about every branch of mathematics that you can think of can be "spoken" abstractly with the language of category theory. While a thorough treatment of the subject is way beyond our scope, I will present two key ingredients that permeate much of algebraic topology (and indeed mathematics in general). These are the concepts of a "category" and that of "maps" between different categories. I will also present some heavy machinery here that will most likely make you cry and beg for mercy, but don't worry, that madness won't be required on subsequent sections. I merely present them to give you a "bird's eye view" of the wonderful subject of category theory, and to encourage you to view mathematics from a "higher" point of abstraction.

**Definition 91.** *A **category** $\mathcal{C}$ consists of*

- *a class of **objects**, denoted $\mathrm{Ob}(\mathcal{C})$.*

- *given two objects $x, y \in \mathcal{C}$, a set $\mathrm{Hom}(x, y)$ of **morphisms**. Generalizing from the categories where $\mathrm{Hom}(x, y)$ is a set of functions (which is not always the case, as you will see shortly!), we usually denote $f \in \mathrm{Hom}(x, y)$ by $f \colon x \to y$. Morphisms satisfy the following properties:*

    - *given morphisms $f \colon x \to y$ and $g \colon y \to z$, we can compose them and obtain $g \circ f \colon x \to$*

*z. When there is no possibility of confusion $g \circ f$ is abbreviated $gf$.*



    **–** *for any $x \in C$, there is an **identity** morphism $1_x \colon x \to x$ such that, for any $f \colon x \to y$, we have $f 1_x = f = 1_y f$.*

As I alluded to in the definition, morphisms are more general than mappings. As you will see in Example 46 below, there are morphisms that you are already quite familiar with, but there are some that will make you pull your hair off. I will present one such outrageous example here because it comes from a topic in algebraic topology that I am personally quite fond of, and I would like to introduce it to you at this point. I am not expecting you to understand this material at this stage, since it relies on some heavy machinery that we haven't even covered and that is way beyond the scope of this humble little book, but nonetheless it will show you that the concept of morphisms is not as simple as you may have thought upon brief inspection of Definition 91. Ready to be mindblown again? Here we go!:

**Definition 92.** *An $(n+1)$-dimensional **cobordism** is a quintuple $(W; M, N, \iota_M, \iota_N)$ consisting of an $(n+1)$-dimensional compact smooth manifold with boundary $W$, closed $n$-manifolds $M$ and $N$, and embeddings $\iota_M \colon M \hookrightarrow \partial W$ and $\iota_N \colon N \hookrightarrow \partial W$, with disjoint images such that*

$$\partial W = \iota_M(M) \amalg \iota_N(N).$$

*The quintuple is usually abbreviated to $(W; M, N)$, omitting any mention to the embeddings. The $n$-manifolds $M$ and $N$ are said to be **cobordant** if such a cobordism exists. All manifolds cobordant to a fixed given manifold $M$ form the cobordism class of $M$.*

Whew, that was intense, wasn't it?? Again, don't worry one bit if you didn't understand a single word on this definition; everything will seem much less daunting after you get through Chapter 3, where you will deal with these so-called "manifolds" and with some of the heavy, well-oiled mathematical machinery associated to them. For now, my purpose is not to give you brain damage, but rather to spark your curiosity for what's ahead and

Jackie Chan's reaction to Definition 92.

also provide an example that showcases the complexity of some not-so-familiar categories. Without further ado, let's give a brief example:

**Example 46.** *Examples of categories are:*

- Set, *where the objects are sets and the morphisms are functions (this is a category that you've been dealing with since kindergarten!).*

- nCob, *where the objects are* $(n-1)$*-dimensional oriented compact manifolds, and morphisms are (cobordism classes of) n-dimensional cobordisms. Here's a pretty picture of cobordisms with* $n = 2$:

$$S \amalg S$$
$$\downarrow M^* \amalg 1_S$$
$$S \amalg S \amalg S$$
$$\downarrow 1_S \amalg M$$
$$S \amalg S$$



*The identity morphism in* 2Cob *is given by the (cobordism class of the) cylinder:*

$$S$$
$$\downarrow [0,1] \times S$$
$$S$$



- Vect, *where objects are (finite-dimensional) vector spaces, and morphisms are linear operators.*

- Hilb, *where objects are (finite-dimensional) Hilbert spaces, and morphisms are linear operators.*

- Top, *where the objects are topological spaces and the morphisms are continuous maps.*

- Ring, *where the objects are rings and the morphisms are ring homomorphisms.*

- Grp, *where the objects are groups and the morphisms are group homomorphisms.*

- Top$_*$, *where the objects are **pointed** topological spaces and the morphisms are **pointed** continuous maps.*

These last two categories are the ones that we're are really after. In particular we will want to "send" objects from the category of pointed topological spaces (a concept that will be discussed later on) to objects that lie in the category of groups. In other words, we would like to have at our disposal a "map" Top$_* \to$ Grp; this is precisely what our next definition will provide for us:

**Definition 93.** *Given categories $\mathcal{C}$ and $\mathcal{D}$, a **functor** $F \colon \mathcal{C} \to \mathcal{D}$ consists of:*

- *A map sending any object $x \in \mathcal{C}$ to an object $F(x) \in \mathcal{D}$.*

- *For any pair of objects $x$ and $y$, a map sending morphisms $f \colon x \to y$ to morphisms $F(f) \colon F(x) \to F(y)$, such that these laws hold:*

    - *for any object $x \in \mathcal{C}$, we have $F(1_x) = 1_{F(x)}$.*
    - *for any pair of morphisms $f \colon x \to y$ and $g \colon y \to z$, we have $F(gf) = F(g)F(f)$.*

*In short: F sends objects to objects, morphisms to morphisms, and preserves sources, targets, identities, and composition.*

**Definition 94.** *A functor $F \colon C \to D$ is called **faithful** if for each pair of objects $X, Y \in C$, the map $F_{X,Y} \colon \mathrm{Hom}_C(X, Y) \to \mathrm{Hom}_D(F(X), F(Y))$ is injective. F is called **full** if the maps $F_{X,Y} \colon \mathrm{Hom}_C(X, Y) \to \mathrm{Hom}_D(F(X), F(Y))$ are all surjective. If C is a subcategory of D, then the inclusion functor is always faithful. If it is full, C is called a **full subcategory**.*

**Definition 95.** *A functor is called **essentially surjective** if every object in D is isomorphic to an image under F of an object of C. A functor is called an **equivalence** if it is faithful, full, and essentially surjective.*

If there is one concept that you should absolutely take away from this subsection is that of a functor. These dudes are literally everywhere in mathematics, and you encounter them more often than you think! Don't worry if all the categorical nonsense you've seen here has given you a massive headache; I've been dealing with this crazy stuff for a while and it still makes my head spin more than a pulsar! You can safely skip the rest of this subsection, since there will be nothing relevant to our subsequent topics...or you can stay for a little more ...☺

**Definition 96.** *We say that a category $C$ has **adjoints** or **duals for morphisms** or that is a $*$-category if there is a contravariant functor $*: C \to C$ which takes objects to themselves and such that $*^2 = 1$ (the identity functor). For any object $x$ or morphism $f$, the dual is denoted $*(x) = x^*$ or $*(f) = f^*$.*

Spelling out the definition, the contravariant functor $*$ has to satisfy the following properties:

- $x^* = x$ for any $x \in C$,

- for any $f: x \to y$ there is a morphism $f^*: y \to x$ (this is what "contravariant" means),

- for any $x \in C$, $(1_x)^* = 1_{x^*} = 1_x$,

- for any morphisms $f: x \to y$ and $g: y \to z$, we have $(gf)^* = f^* g^*$, and

- $(f^*)^* = f$ for any morphism $f$.

**Example 47.** *Examples of $*$-categories are:*

- nCob, *where $M^*$ is obtained by exchanging the roles of input and output. If the cobordism is imbedded, this can be represented as reflection along the "time" direction.*



- *any **groupoid** (a category where every morphism is invertible), as then the inverse has the properties required of $*$.*

- Hilb, *where the adjoint $T^*$ of a linear operator $T\colon \mathcal{H} \to \mathcal{H}'$ is defined by $\langle T^*\phi \mid \psi \rangle_{\mathcal{H}} = \langle \phi \mid T\psi \rangle_{\mathcal{H}'}$. If this notation looks weird to you it's because it \*is\* weird; what do you expect? It was invented by physicists! (I'm joking, please lay down the pitchforks!). This is the so called **Bra-Ket notation** or **Dirac notation**. Here's a link to a short-and-sweet 3-minute video explaining this madness.*

**Definition 97.** *We say that $F\colon \mathcal{C} \to \mathcal{D}$ is a ∗-**functor** if given $f\colon x \to y$, we have $F(f^*) = F(f)^*\colon F(x) \to F(y)$.*

**Definition 98.** *A category $\mathcal{C}$ is **monoidal** if it is equipped with an operation $\otimes$ with the following properties:*

- *for any $x, y \in \mathcal{C}$, there is an object $x \otimes y \in \mathcal{C}$;*

- *for any $f\colon x \to x'$ and $g\colon y \to y'$, there is a morphism $f \otimes g\colon x \otimes y \to x' \otimes y'$.*

- *for any objects $x, y, z \in \mathcal{C}$ there is an isomorphism $a_{xyz}\colon (x \otimes y) \otimes z \to x \otimes (y \otimes z)$ called the **associator** and satisfying the **pentagon identity** :*

$$
\begin{array}{ccccc}
 & & (w\otimes x)\otimes(y\otimes z) & & \\
 & \nearrow^{a_{(wx)yz}} & & \searrow^{a_{wx(yz)}} & \\
\big((w\otimes x)\otimes y\big)\otimes z & & & & w\otimes\big(x\otimes(y\otimes z)\big) \\
\downarrow^{a_{wxy}\otimes 1_z} & & & & \nearrow^{1_w\otimes a_{xyz}} \\
\big(w\otimes(x\otimes y)\big)\otimes z & \xrightarrow{\;\;a_{w(xy)z}\;\;} & w\otimes\big((x\otimes y)\otimes z\big) & &
\end{array}
$$

- *there is an object $1$ such that, for any object $x \in C$, there are isomorphisms $l_x\colon 1 \otimes x \to x$ and*

$r_x \colon x \otimes 1 \to x$ called **units** satisfying the other identity:

$$
\begin{array}{ccccc}
 & & x \otimes 1 & & \\
 & \overset{l_x \otimes 1_1}{\nearrow} & & \overset{r_x}{\searrow} & \\
(1 \otimes x) \otimes 1 & & & & x \\
 & \underset{a_{1x1}}{\searrow} & & \underset{l_x}{\nearrow} & \\
 & & 1 \otimes (x \otimes 1) \xrightarrow{\;1_1 \otimes r_x\;} 1 \otimes x & &
\end{array}
$$

- *finally, given $f \colon x \to y$, $g \colon y \to z$, $f' \colon x' \to y'$ and $g' \colon y' \to z'$, we require that $(g \otimes g')(f \otimes f') = (gf) \otimes (g'f')$, which just says that the following diagram is unambiguous:*



**Example 48.** *Examples of monoidal categories are*

- Group*: objects are groups, morphisms are group homomorphisms and $\otimes$ is the direct product of groups.*

- nCob*: the $\otimes$, both for objects and for morphisms, is the disjoint union of manifolds, i.e., $\otimes = \amalg$.*

- Vect *or* Hilb*: the $\otimes$ is the tensor product. This is how in quantum mechanics two things are "put together."*

One last badass definition to close out this subsection:

**Definition 99.** *A **topological quantum field theory** (or **TQFT** for short) is a monoidal functor from a category of cobordisms to a category of vector spaces. That is, it is a functor whose value on a disjoint union of manifolds is equivalent to the tensor product of its values on each of the constituent manifolds.*

To illustrate our definition, consider a functor $Z \colon 2\mathrm{Cob} \to \mathrm{Vect}_{\Bbbk}$ (by this I mean vector spaces with underlying field $\Bbbk$. The following diagram brings to life our definition:

$$
\begin{array}{ccc}
Z \colon & 2\mathrm{Cob} & \longrightarrow & \mathrm{Vect}_{\Bbbk} \\
& & \mapsto & V \\
& & \mapsto & V \otimes V \to V \\
& & \mapsto & \Bbbk \to V \\
& & \mapsto & V \to \Bbbk
\end{array}
$$

For an even more elaborate construction, consider the figure you saw earlier on Example 46; this time we add the target of the TQFT (i.e., the vector spaces) to the figure so you can clearly see Definition 99 in all its glory.

$$
\begin{array}{ccc}
S \amalg S & & V \otimes V \\
\downarrow {\scriptstyle M^* \amalg 1_S} & & \downarrow \\
S \amalg S \amalg S & & V \otimes V \otimes V \\
\downarrow {\scriptstyle 1_S \amalg M} & & \downarrow \\
S \amalg S & & V \otimes V
\end{array}
$$

### 2.3.2 Basics of Homotopy

**Definition 100.** *If $\pi \colon X \to Y$ is a map, a subset $U \subseteq X$ is said to be **saturated with respect to** $\pi$ if $U$ is the entire preimage of its image under $\pi$; i.e., if $U = \pi^{-1}(\pi(U))$. Given $y \in Y$, the **fiber of $\pi$ over $y$** is the set $\pi^{-1}(y)$. (Thus, a subset of $X$ is saturated if and only if it is a union of fibers).*

**Definition 101.** *If $X$ and $Y$ are topological spaces, a map $F\colon X \to Y$ (continuous or not) is said to be **proper** if for every compact set $K \subseteq Y$, the preimage $F^{-1}(K)$ is compact as well.*

Here are some useful sufficient conditions for a map to be proper:

**Proposition 30 (Sufficient Conditions for Properness).** *Suppose $X$ and $Y$ are topological spaces, and $F\colon X \to Y$ is a continuous map.*

  *a)* *If $X$ is compact and $Y$ is Hausdorff, then $F$ is proper.*

  *b)* *If $F$ is a closed map with compact fibers, then $F$ is proper.*

  *c)* *If $F$ is a topological embedding with closed image, then $F$ is proper.*

  *d)* *If $Y$ is Hausdorff and $F$ has a continuous left inverse (i.e., a continuous map $G\colon Y \to X$ such that $G \circ F = \mathrm{Id}_X$), then $F$ is proper.*

  *e)* *If $F$ is proper and $A \subseteq X$ is a subset that is saturated with respect to $F$, then $F|_A\colon A \to F(A)$ is proper.*

**Definition 102.** *A **deformation retraction** of a space $X$ onto a subspace $A$ is a family of maps $f_t\colon X \to X$, for $t \in I$, such that $f_0(X) = \mathrm{Id}$ (the identity map), $f_1(X) = A$, and $f_t|_A = \mathrm{Id}$ for all $t$. The family $f_t$ should be continuous in the sense that the associated map $X \times I \to X$ given by $(x, t) \mapsto f_t(x)$, is continuous.*

**Example 49.** *We now construct an explicit deformation retraction of the torus with one point deleted onto a graph consisting of two circles intersecting in a point, namely, longitude and meridian circles of the torus. Using the CW complex construction of the torus (see this on [Hatcher, 2001, Pg 5]), we have the map denoted by the small arrows:*

*To prove this map is indeed a deformation retraction, we use the identification of the unit square with the unit disc (in this case the boundary of the circle is divided into four arcs with the labeling scheme $aba^{-1}b^{-1}$), and using polar coordinates we can let $\tilde{F}((r,\theta),t) = (r + t(1 - r),\theta)$. Then, let $F = q \circ \tilde{F}$ where $q$ is the quotient map shown as the last arrow above. $F$ is continuous as $\tilde{F}$ is continuous in each coordinate, and then $q$ is continuous. $F$ is also such that $F((r,\theta),0) = (r,\theta)$, $F((r,\theta),1) = (1,\theta)$, and $F((1,\theta),t) = (1,\theta)$. Thus, $F$ is a deformation retraction. The two circles in the last diagram are the longitude and meridian circles of the torus by construction, since on the third diagram of the figure we identify the sides marked a to get the meridian circle and then the sides marked b to get the longitude circle.* ◉

**Definition 103.** *For a map $f\colon X \to Y$, the **mapping cylinder** $M_f$ is the quotient space of the disjoint union $(X \times I) \amalg Y$ obtained by identifying each $(x,1) \in X \times I$ with $f(x) \in Y$.*



**Definition 104.** *A **homotopy** is any family of maps $f_t\colon X \to Y$, for $t \in I$, such that the associated map $F\colon X \times I \to Y$ given by $F(x,t) = f_t(x)$ is continuous. One says that two maps $f_0, f_1\colon X \to Y$ are **homotopic** if there exists a homotopy $f_t$ connecting them, in which case we write $f_0 \simeq f_1$.*

In these terms, a deformation retraction of $X$ onto a subspace $A$ is a homotopy from the identity map of $X$ to a **retraction** of $X$ onto $A$ (a map $r\colon X \to X$ such that $r(X) = A$ and $r|_A = \mathrm{Id}$. Equivalently, we may regard a retraction as a map $X \to A$ restricting to the identity on the subspace $A \subset X$; the subspace $A \subset X$ in this case is called a **retract** of $X$.) From a more formal viewpoint a retraction is a map $r\colon X \to X$ with $r^2 = r$, since this equation says exactly that $r$ is the identity on its image. Retractions are the topological analogs of projection operators in other parts of mathematics.

**Remark**: A homotopy $f_t\colon X \to X$ that gives a deformation retraction of $X$ onto a subspace $A$ has the property that $f_t|_A = \mathrm{Id}$ for all $t$. In general, a homotopy $f_t\colon X \to Y$ whose

restriction to a subspace $A \subset X$ is independent of $t$ is called a **homotopy relative to** $A$ (or more concisely, a homotopy rel $A$). Thus, a deformation retraction of $X$ onto $A$ is a homotopy rel $A$ from the identity map of $X$ to a retraction of $X$ onto $A$.

If a space $X$ deformation retracts onto a subspace $A$ via $f_t \colon X \to X$, then if $r \colon X \to A$ denotes the resulting retraction and $\iota \colon A \to X$ the inclusion, we have $r\iota = \mathrm{Id}$ and $\iota r \simeq \mathrm{Id}$, the latter homotopy being given by $f_t$. Generalizing this situation, a map $f \colon X \to Y$ is called a **homotopy equivalence** if there is a map $g \colon Y \to X$ such that $fg \simeq \mathrm{Id}$ and $gf \simeq \mathrm{Id}$. The spaces $X$ and $Y$ are said to be **homotopy equivalent** or to have the same **homotopy type**, which we denote by $X \simeq Y$. It is true in general that two spaces $X$ and $Y$ are homotopy equivalent if and only if there exists a third space $Z$ containing both $X$ and $Y$ as deformation retracts.

**Definition 105.** *A space having the homotopy type of a point is called **contractible**. This amounts to requiring that the identity map of the space be **nullhomotopic**, that is, homotopic to a constant map.*

In general, this is slightly weaker than saying the space deformation retracts to a point; see the exercises at the end of Chapter 1 of [Hatcher, 2001]., for an example distinguishing these two notions.

**Example 50.** *We now show that a retract of a contractible space is contractible. Saying that a space $X$ is contractible is equivalent to saying that $\mathrm{Id}_X$ is homotopic to a constant map mapping to some point $x_0 \in X$. Let $F(x,t)$ be this homotopy, and let $A \subset X$ be our retract, i.e., there exists a map $r \colon X \to A$ such that $r|_A = \mathrm{Id}_A$. Then, $(F \circ r)|_A$ is a homotopy between $\mathrm{Id}_A$ and the constant map mapping to $r(x_0)$, for it is a composition of continuous maps hence itself continuous, and since $r(F(x,0))|_A = (\mathrm{Id}_X \circ r)|_A = r|_A = \mathrm{Id}_A$ and $r(F(x,1))|_A = r(x_0)$.* 🌐

**Example 51.** *In this example we show that a retract of a Hausdorff space is closed. Let $A \subset X$ be a retract of X. Consider a point $x \in \partial A$. If $x \notin A$, then $r(x) \neq x$ and there are disjoint neighborhoods $U$ of $x$ and $V$ of $r(x)$. Then by continuity of $r$, there must be a neighborhood $W \subseteq U$ of $x$, so that $r(W) \subseteq V$. However, $r(W \cap A)$ is a nonempty subset of $W$,[2] so it cannot be in $V$. Hence we must have $r(x) = x \in A$, and hence $A$ is closed.*

✦──────────────── ☞

---

[2] $W \cap A$ is nonempty because, by assumption, $x \in \partial A$.

Alternatively, we could have used the fact that for any two maps $f, g \colon X \to A$, the so-called **equalizer** $\{x \in X \mid f(x) = g(x)\}$ is a closed subspace of $X$ if $A$ is Hausdorff. This follows from the diagonal $\Delta_A$ being closed in $A \times A$ and $(f, g) \colon X \to A \times A$ being continuous. In this particular case we just take $f = \mathrm{Id}_A$ and $g = r$; the equalizer is then the retract $A$. ✿

### 2.3.3 Covering Maps

**Definition 106.** *Suppose $E$ and $X$ are topological spaces. A map $\pi \colon E \to X$ is called a **covering map** if $E$ and $X$ are connected and locally path-connected, $\pi$ is surjective and continuous, and each point $p \in X$ has a neighborhood $U$ that is **evenly covered by** $\pi$, meaning that each component of $\pi^{-1}(U)$ is mapped homeomorphically onto $U$ by $\pi$. In this case, $X$ is called the **base of the covering**, and $E$ is called a **covering space of** $X$. If $U$ is an evenly covered subset of $X$, the components of $\pi^{-1}(U)$ are called the **sheets of the covering over** $U$.*



Figure 2.11: An evenly covered neighborhood of $x$. (In this figure $\pi = q$)

**Example 52.** *Let $E$ be the interval $(0, 2) \subset \mathbb{R}$, and define $f \colon E \to \mathbb{S}^1$ by $f(x) = e^{2\pi i x}$ (see Figure 2.12). Then $f$ is a local homeomorphism (because it is the restriction of the covering map $\epsilon \colon \mathbb{R} \to \mathbb{S}^1$ given by $\epsilon(x) = e^{2\pi i x}$), and is clearly surjective. However, $f$ is not a covering map because the point $1 \in \mathbb{S}^1$ has no evenly covered neighborhood.* ✿

Figure 2.12: A surjective local homeomorphism that is not a covering map.

**Definition 107.** *If $\pi\colon E \to X$ is a covering map and $F\colon B \to X$ is a continuous map, a **lift of F** is a continuous map $\widetilde{F}\colon B \to E$ such that $\pi \circ \widetilde{F} = F$:*



**Proposition 31 (Lifting Properties of Covering Maps).** *Suppose $\pi\colon E \to X$ is a covering map.*

   *a)* UNIQUE LIFTING PROPERTY: *If $B$ is a connected space and $F\colon B \to X$ is a continuous map, then any two lifts of $F$ that agree at one point are identical.*

   *b)* PATH LIFTING PROPERTY: *If $f\colon I \to X$ is a path, then for any point $e \in E$ such that $\pi(e) = f(0)$, there exists a unique lift $\widetilde{f}\colon I \to E$ of $f$ such that $\widetilde{f}(0) = e$.*

   *c)* MONODROMY THEOREM: *If $f,g\colon I \to X$ are path-homotopic paths and $\widetilde{f}_e, \widetilde{g}_e\colon I \to E$ are their lifts starting at the same point $e \in E$, then $\widetilde{f}_e$ and $\widetilde{g}_e$ are path-homotopic and $\widetilde{f}_e(1) = \widetilde{g}_e(1)$.*

**Proposition 32 (Lifting Criterion).** *Suppose $\pi\colon E \to X$ is a covering map, $Y$ is a connected and locally path-connected space, and $F\colon Y \to X$ is a continuous map. Let $y \in Y$ and $e \in E$ be such that $\pi(e) = F(y)$. Then there exists a lift $\widetilde{F}\colon Y \to E$ of $F$ satisfying $\widetilde{F}(y) = e$ if and only if $F_*\left(\pi_1(Y,y)\right) \subseteq \pi_*\left(\pi_1(E,e)\right)$.*

**Proposition 33 (Coverings of Simply Connected Spaces).** *If $X$ is a simply connected space, then every covering map $\pi\colon E \to X$ is a homeomorphism.*

**Definition 108.** *A topological space is said to be **locally simply connected** if it admits a basis of simply connected open subsets.*

**Proposition 34 (Existence of a Universal Covering Space).** *If $X$ is a connected and locally simply connected topological space, there exists a simply connected topological space $\widetilde{X}$ and a covering map $\pi\colon \widetilde{X} \to X$. If $\widehat{\pi}\colon \widehat{X} \to X$ is any other simply connected covering of $X$, then there is a homeomorphism $\varphi\colon \widetilde{X} \to \widehat{X}$ such that $\widehat{\pi} \circ \varphi = \pi$.*

**Definition 109.** *The simply connected covering space $\widetilde{X}$ whose existence and uniqueness (up to homeomorphism) are guaranteed by this last proposition is called the **universal covering space of $X$.***

# Chapter 3

# Smooth Manifolds

D IFFERENTIAL GEOMETRY is a beautiful language in which much of modern mathematics and physics is spoken. It is my personal favorite subject in all of mathematics, and it is my hope that I can properly convey my boundless enthusiasm for it in these notes.

There are different approaches to the study of differential geometry. Some authors approach the subject with a treatment of curves and surfaces in $\mathbb{R}^3$, which provide much of the motivation and intuition for the general theory. Others choose to go straight into the more general point of view of abstract manifolds, which are generalizations of curves and surfaces to arbitrarily many dimensions that provide the mathematical context for understanding "space" in all of its manifestations. These notes follow the latter approach. To study the geometry of curves and surfaces, I highly recommend books such as *Differential Geometry of Curves and Surfaces* by Manfredo P. do Carmo and *Elementary Differential Geometry* by Andrew Pressley.

Be warned that if I claimed earlier that real analysis was calculus on steroids, I claim now that differential geometry is multivariate calculus on meth (and I wouldn't have it any other way!). Make sure that you already have a solid acquaintance with undergraduate linear algebra, real analysis, multivariate calculus, and topology before reading this chapter. Only then you will be able to truly appreciate how marvelous this subject really is.

## 3.1   Calculus Preliminaries

If $f\colon A \to \mathbb{R}$ is bounded, the extent to which $f$ fails to be continuous at $a \in A$ can be measured in a precise way: For $\delta > 0$, let

$$M(a,f,\delta) = \sup\{f(x) \mid x \in A, |x - a| < \delta\},$$
$$m(a,f,\delta) = \inf\{f(x) \mid x \in A, |x - a| < \delta\}.$$

**Definition 110.** *The **oscillation** $o(f,a)$ of $f$ at $a$ is defined by $o(f,a) = \lim_{\delta \to 0}[M(a,f,\delta) - m(a,f,\delta)]$.*

This limit always exists, since $M(a,f,\delta) - m(a,f,\delta)$ decreases as $\delta$ decreases. There are two important facts about $o(f,a)$, given in the following two theorems:

**Theorem 89.** *The bounded function $f$ is continuous at $a$ if and only if $o(f,a) = 0$.*

**Theorem 90.** *Let $A \subset \mathbb{R}^n$ be closed. If $f\colon A \to \mathbb{R}$ is any bounded function, and $\varepsilon > 0$, then $\{x \in A \mid o(f,x) \geq \varepsilon\}$ is closed.*

### 3.1.1   Differentiation

**Definition 111.** *A function $f : \mathbb{R}^n \to \mathbb{R}^m$ is **differentiable** at $a \in \mathbb{R}^n$ if there exists a linear transformation $\lambda\colon \mathbb{R}^n \to \mathbb{R}^m$ such that*

$$\lim_{h \to 0} \frac{\|f(a + h) - f(a) - \lambda(h)\|}{\|h\|} = 0. \tag{3.1}$$

**Remark**: Note that $h$ is a point of $\mathbb{R}^n$ while $f(a + h) - f(a) - \lambda(h)$ is a point of $\mathbb{R}^m$, hence the norm signs on (3.1) are essential. The linear transformation $\lambda$ is usually denoted $Df(a)$ and it is called the **derivative** of $f$ at $a$. Also, the matrix associated with such linear transformation is called the ***Jacobian matrix*** of $f$ at $a$.

The justification for the phrase "the linear transformation $\lambda$" is given by the following theorem:

**Theorem 91.** *If $f : \mathbb{R}^n \to \mathbb{R}^m$ is differentiable at $a \in \mathbb{R}^n$, there is a unique linear transformation $\lambda : \mathbb{R}^n \to \mathbb{R}^m$ such that equation (3.1) is satisfied.*

*Proof.* Suppose $\mu : \mathbb{R}^n \to \mathbb{R}^m$ satisfies

$$\lim_{h \to 0} \frac{\|f(a+h) - f(a) - \mu(h)\|}{\|h\|} = 0.$$

If $d(h) = f(a+h) - f(a)$, then

$$
\begin{aligned}
\lim_{h \to 0} \frac{\|\lambda(h) - \mu(h)\|}{\|h\|} &= \lim_{h \to 0} \frac{\|\lambda(h) - d(h) + d(h) - \mu(h)\|}{\|h\|} \\
&\leq \lim_{h \to 0} \frac{\|\lambda(h) - d(h)\|}{\|h\|} + \lim_{h \to 0} \frac{\|d(h) - \mu(h)\|}{\|h\|} \\
&= 0.
\end{aligned}
$$

If $x \in \mathbb{R}^n$, then $tx \to 0$ as $t \to 0$. Hence for $x \neq 0$ we have

$$0 = \lim_{t \to 0} \frac{\|\lambda(tx) - \mu(tx)\|}{\|tx\|} = \frac{\|\lambda(x) - \mu(x)\|}{\|x\|}.$$

Therefore $\lambda(x) = \mu(x)$, and this concludes our proof. $\qquad\square$

**Definition 112.** *Let $f \colon \mathbb{R}^n \to \mathbb{R}$ and $a \in \mathbb{R}$. If the limit*

$$\lim_{h \to 0} \frac{f(a^1, \ldots, a^i + h, \ldots, a^n) - f(a^1, \ldots, a^n)}{h}$$

*exists, then it is called the $i^{th}$ **partial derivative** of $f$ at $a$, and it is denoted by $\partial_i f(a)$.*

**Theorem 92.** *Let $A \subset \mathbb{R}^n$. If an extremum of $f \colon A \to \mathbb{R}$ occurs at a point $a$ in the interior of $A$ and $\partial_i f(a)$ exists, then $\partial_i f(a) = 0$.*

*Proof.* Let $g_i(x) = f(a^1, \ldots, x, \ldots, a^n)$. Clearly $g_i$ has a maximum (or minimum) at $a^i$, and $g_i$ is defined in an open interval containing $a^i$. Hence $0 = g_i'(a^i) = \partial_i f(a)$.     $\square$

**Theorem 93.** *If $f : \mathbb{R}^n \to \mathbb{R}^m$ is differentiable at $a$, then $\partial_j f^i(a)$ exists for $1 \le i \le m$, $1 \le j \le n$, and $f'(a)$ is the $m \times n$ matrix $(\partial_j f^i(a))$.*

There are several examples in the problems to show that the converse of the above theorem is false. It is true, however, if one hypothesis is added, as shown by the following theorem:

**Theorem 94.** *If $f : \mathbb{R}^n \to \mathbb{R}^m$, then $Df(a)$ exists if all $\partial_j f^i(x)$ exist in an open set containing $a$ and if each function $\partial_j f^i$ is continuous at $a$. Such a function $f$ is said to be* **continuously differentiable** *at $a$.*

**Lemma 23.** *Let $A \subset \mathbb{R}^n$ be a rectangle and let $f : A \to \mathbb{R}^n$ be continuously differentiable. If there is a number $M$ such that $\|\partial_j f^i(x)\| \le M$ for all $x$ in the interior of $A$, then*

$$\|f(y) - f(x)\| \le n^2 M \|y - x\| \qquad \forall\ x, y \in A.$$

*Proof.* We have that

$$f^i(y) - f^i(x) = \sum_{j=1}^n [f^i(y^1, \ldots, y^j, x^{j+1}, \ldots, x^n) - f^i(y^1, \ldots, y^{j-1}, x^j, \ldots, x^n)].$$

Now, applying the mean value theorem, there exists some $z_{ij}$ such that

$$f^i(y^1, \ldots, y^j, x^{j+1}, \ldots, x^n) - f^i(y^1, \ldots, y^{j-1}, x^j, \ldots, x^n) = (y^j - x^j) \cdot \partial_j f^i(z_{ij})$$
$$\le M \cdot \|y^j - x^j\|.$$

Thus,

$$\|f^i(y) - f^i(x)\| \le \sum_{j=1}^n \|y^j - x^j\| \cdot M \le nM \|y - x\| \qquad \text{(Since each } \|y^j - x^j\| \le \|y - x\|\text{).}$$

Finally,

$$\|f(y) - f(x)\| \le \sum_{j=1}^n \|f^i(y) - f^i(x)\| \le n^2 M \cdot \|y - x\|. \qquad \square$$

**Theorem 95 (Inverse Function Theorem).** *Suppose that $f\colon \mathbb{R}^n \to \mathbb{R}^n$ is continuously differentiable in an open set containing $a$, and $\det f'(a) \neq 0$. Then there is an open set $V$ containing $a$ and an open set $W$ containing $f(a)$ such that $f\colon V \to W$ has a continuous inverse $f^{-1}\colon W \to V$ which is differentiable and for all $y \in W$ satisfies*

$$(f^{-1})'(y) = \frac{1}{f'(f^{-1}(y))}.$$

*Proof.* See [Spivak, 1971, p. 35]                                             $\square$

It should be noted that an inverse function $f^{-1}$ may exist even if $\det f'(a) = 0$. For example, if $f\colon \mathbb{R} \to \mathbb{R}$ is defined by $f(x) = x^3$, then $f'(0) = 0$ but $f$ has the inverse function $f^{-1}(x) = \sqrt[3]{x}$. One thing is certain however: if $\det f'(a) = 0$, then $f^{-1}$ cannot be differentiable at $f(a)$. To prove this, note that $f \circ f^{-1}(x) = x$. If $f^{-1}$ were differentiable at $f(a)$, then the chain rule would give

$$f'(a) \cdot (f^{-1})'(f(a)) = \mathrm{Id}$$
$$\implies \det f'(a) \cdot \det(f^{-1})'(f(a)) = 1,$$

contradicting the assumption that $\det f'(a) = 0$.

Now we want to ask ourselves the following question:

If $f\colon \mathbb{R}^n \times \mathbb{R} \to \mathbb{R}$ with $f(a^1, \ldots, a^n, b) = 0$, when can we find, for each $(x^1, \ldots, x^n)$ near $(a^1, \ldots, a^n)$, a unique $y$ near $b$ such that $f(x^1, \ldots, x^n, y) = 0$?

Even more generally, we can ask about the possibility of solving $m$ equations, depending upon parameters $x^1, \ldots, x^n$, in $m$ unknowns: If

$$f_i\colon \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R} \qquad \text{for } i = 1, \ldots, m$$

and

$$f_i(a^1, \ldots, a^n, b^1, \ldots, b^m) = 0 \qquad \text{for } i = 1, \ldots, m,$$

when can we find, for each $(x^1, \ldots, x^n)$ near $(a^1, \ldots, a^n)$ a unique $(y^1, \ldots, y^m)$ near $(b^1, \ldots, b^m)$ which satisfies $f_i(x^1, \ldots, x^n, y^1, \ldots, y^m) = 0$? The answer is provided by the following theorem:

**Theorem 96 (Implicit Function Theorem).** *Suppose that* $f\colon \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}^m$ *is continuously differentiable in an open set containing* $(a, b)$*, and* $f(a, b) = 0$*. Let M be the* $m \times m$ *matrix*

$$\left( \partial_{n+j} f^i(a, b) \right) \qquad \text{for } 1 \le i, j \le m.$$

*If* $\det M \neq 0$*, there is an open set* $A \subset \mathbb{R}^n$ *containing a and an open set* $B \subset \mathbb{R}^m$ *containing b, with the following property: for* $x \in A$ *there is a unique* $g(x) \in B$ *such that* $f(x, g(x)) = 0$*. The function g is differentiable.*

*Proof.* Define $F\colon \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}^n \times \mathbb{R}^m$ by $F(x, y) = (x, f(x, y))$. Then $\det F'(a, b) = \det M \neq 0$. By the *Inverse Function Theorem* there is an open set $W \subset \mathbb{R}^n \times \mathbb{R}^m$ containing $F(a, b) = (a, 0)$ and an open set in $\mathbb{R}^n \times \mathbb{R}^m$ containing $(a, b)$, which we may take to be of the form $A \times B$, such that $F\colon A \times B \to W$ has a differentiable inverse $h\colon W \to A \times B$. Clearly $h$ is of the form $h(x, y) = (x, k(x, y))$ for some differentiable function $k$ (since $F$ is of this form). Now let $\pi_2\colon \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}^m$ be the projection on the second component defined by $\pi_2(x, y) = y$; then $\pi_2 \circ F = f$. Therefore

$$
\begin{aligned}
f(x, k(x, y)) &= f \circ h(x, y) = (\pi_2 \circ F) \circ h(x, y) \\
&= \pi_2 \circ (F \circ h)(x, y) = \pi_2(x, y) = y.
\end{aligned}
$$

Thus, $f(x, k(x, 0)) = 0$; in other words we can define $g(x) = k(x, 0)$. $\qquad \square$

### 3.1.2 Integration

**Definition 113.** *The **volume** of a closed k-cube* $[a_1, b_1] \times \cdots \times [a_k, b_k]$*, and also of an open k-cube* $(a_1, b_1) \times \cdots \times (a_k, b_k)$*, is given by* $(b_1 - a_1) \cdot \cdots \cdot (b_k - a_k)$*.*

**Definition 114.** *Suppose that A is a rectangle,* $f\colon A \to \mathbb{R}$ *is bounded, and P is a partition of A. For each subrectangle S of the partition, let*

$$
\begin{aligned}
m_S(f) &= \inf\{f(x) : x \in S\}, \\
M_S(f) &= \sup\{f(x) : x \in S\},
\end{aligned}
$$

and let $v(S)$ be the volume of $S$. Then the **lower sum** and **upper sum** of $f$ for $P$ are defined by

$$L(f, P) = \sum_S m_S(f) \cdot v(S) \quad and \quad U(f, P) = \sum_S M_S(f) \cdot v(S),$$

*respectively.*

Clearly, $L(f, P) \le U(f, P)$, and an even stronger assertion (given in Corollary 27 below) is true.

**Lemma 24.** *Suppose the partition $P'$ refines $P$ (that is, each subrectangle of $P'$ is contained in a subrectangle of P). Then*

$$L(f, P) \le L(f, P') \quad and \quad U(f, P') \le U(f, P).$$

**Corollary 27.** *If $P$ and $P'$ are any two partitions, then it is always true that $L(f, P') \le U(f, P)$.*

**Definition 115.** *A function $f: A \to \mathbb{R}$ is said to be **integrable** (in the Riemann sense) on the rectangle $A$ if $f$ is bounded and if*

$$\sup\{L(f, P)\} = \inf\{U(f, P)\}.$$

**Theorem 97.** *A bounded function $f: A \to \mathbb{R}$ is Riemann integrable if and only if, for every $\varepsilon > 0$, there is a partition $P$ of $A$ such that $U(f, P) - L(f, P) < \varepsilon$.*

**Example 53.** *a) Here's an example of a Riemann integrable function. Let $f: A \to \mathbb{R}$ be a constant function, $f(x) = c$. Then for any partition $P$ and subrectangle $S$ we have $m_S(f) = M_S(f) = c$, so that*

$$L(f, P) = U(f, P) = \sum_S c \cdot v(S) = c \cdot v(A).$$

*Hence we have $\int_A f = c \cdot v(A)$.*

*b) Here's an example of a nonintegrable function (in the Riemann sense). Let $f: [0, 1] \times [0, 1] \to \mathbb{R}$ be defined by*

$$f(x) = \begin{cases} 0 & if \ x \in \mathbb{Q}, \\ 1 & if \ x \notin \mathbb{Q}. \end{cases}$$

*If P is any partition, then every subrectangle S will contain points $(x, y)$ with x rational, and also points $(x, y)$ with x irrational. Hence, $m_S(f) = 0$ and $M_S(f) = 1$, so*

$$L(f, P) = \sum_S 0 \cdot v(S) = 0,$$

*while*

$$U(f, P) = \sum_S 1 \cdot v(S) = v([0, 1] \times [0, 1]) = 1,$$

*Therefore, f is not Riemann integrable.*

**Lemma 25.** *Let A be a closed rectangle and let $f \colon A \to \mathbb{R}$ be a bounded function such that $o(f, x) < \varepsilon$ for all $x \in A$. Then there is a partition P of A with $U(f, P) - L(f, P) < \varepsilon \cdot v(A)$.*

**Definition 116.** *A subset $A \subset \mathbb{R}^n$ is said to have **measure zero** if for every $\varepsilon > 0$, there is a cover $\{U_1, U_2, \dots\}$ of A by closed (or open) rectangles such that $\sum_{i=1}^{\infty} v(U_i) < \varepsilon$.*

**Definition 117.** *A subset $A \subset \mathbb{R}^n$ is said to have **content zero** if for every $\varepsilon > 0$, there is a finite cover $\{U_1, \dots, U_n\}$ of A by closed (or open) rectangles such that $\sum_{i=1}^{n} v(U_i) < \varepsilon$.*

**Theorem 98.** *If $a < b$, then $[a, b] \subset \mathbb{R}$ does not have content 0. In fact, if $\{U_1, \dots, U_n\}$ is a finite cover of $[a, b]$ by closed intervals, then $\sum_{i=1}^{n} v(U_i) \geq b - a$.*

If $a < b$, it is also true that $[a, b]$ does not have measure 0. This follows from the following theorem:

**Theorem 99.** *If A is compact and has measure 0, then A also has content 0.*

**Remark**: The conclusion of *Theorem 99* is false if A is not compact. For example, let A be the set of rational numbers between 0 and 1; then A has measure 0. Suppose, however, that $\{[a_1, b_1], \dots, [a_n, b_n]\}$ covers A. Then A is contained in the closed set $[a_1, b_1] \bigcup \cdots \bigcup [a_n, b_n]$ and therefore,

$$[0, 1] \subset [a_1, b_1] \bigcup \cdots \bigcup [a_n, b_n].$$

It follows from *Theorem 98* that $\sum_{i=1}^{n}(b_i - a_i) \geq 1$ for any such cover, and consequently A does not have content 0.

**Theorem 100.** *Let $A$ be a closed rectangle and $f \colon A \to \mathbb{R}$ be a bounded function. Let $B = \{x \mid f$ is not continuous at $x\}$. Then $f$ is integrable (in the Riemann sense) if and only if $B$ is a set of measure $0$.*

## 3.2 Introduction to Smooth Manifolds

Before we delve deeper into the subject, let me warn you about the notation we are using. Notation is something that can vary tremendously from author to author all across mathematics, and it is especially true of differential geometry. This is why differential geometry is often referred to as the study of geometric properties that are invariant under change of notation (this is the part where you're supposed to laugh at my terrible jokes).



Math jokes are the best, aren't they?

You can pick up ten different references to study geometry and chances are that they will all use a different notation. That is something that you simply have to deal with, but it isn't the end of the world. As long as you really understand the concepts, the notation won't throw you off much. Just keep an eye out for that and you'll be fine.

### 3.2.1 Some Preliminaries

**Definition 118.** *A topological space $X$ is*

- *connected if there do not exist two disjoint, nonempty, open subsets of $X$ whose union is $X$.*

- *path-connected if every pair of points in $X$ can be joined by a path in $X$.*

- *locally path-connected if $X$ has a basis of path-connected open subsets.*

**Proposition 35.** *If M is a topological manifold, then*

a) *M is locally path-connected.*

b) *M is connected if and only if it is path-connected (In general, connected $\not\Rightarrow$ path-connected, but in manifolds this is always true.)*

c) *the connected components of M are exactly the path-connected components.*

d) *M has countably many components, each of which is open in M and is a connected topological manifold.*

**Definition 119.** *A topological space X is said to be* **locally compact** *if every point has a neighborhood contained in a compact subset of X.*

**Definition 120.** *A subset of a topological space X is said to be* **precompact** *in X if its closure in X is compact.*

**Proposition 36.** *For a Hausdorff space X the following are equivalent:*

a) *X is locally compact.*

b) *Each point of X has a precompact neighborhood.*

c) *X has a basis of precompact open subsets.*

**Lemma 26.** *Every topological manifold has a countable basis of precompact coordinate balls.*

**Proposition 37.** *Every topological manifold is locally compact.*

*Proof.* From Lemma 26 we have that every point in a topological manifold is contained in a precompact coordinate ball.                                                                 $\square$

Let X be a topological space. Then we have the following definitions:

**Definition 121.** *A collection $\{S_i\}$ of subsets of $X$ is said to be **locally finite** if each point of $X$ has a neighborhood that intersects at most finitely many of the sets in $\{S_i\}$.*

**Definition 122.** *Given a cover $\mathcal{U}$ of $X$, another cover $\mathcal{V}$ is called a **refinement** of $\mathcal{U}$ if for each $V \in \mathcal{V}$ there exists some $U \in \mathcal{U}$ such that $V \subseteq U$.*

**Definition 123.** *We say that $X$ is **paracompact** if every open cover of $X$ admits an open, locally finite refinement.*

**Theorem 101.** *Every topological manifold is paracompact. In fact, given a topological manifold $M$, an open cover $\mathcal{U}$ of $M$, and any basis $\mathcal{B}$ for the topology of $M$, there exists a countable, locally finite refinement of $\mathcal{U}$ consisting of elements of $\mathcal{B}$.*

**Theorem 102.** *Let $M$ be a topological manifold. Then its fundamental group $\pi_1(M)$ has countably many elements.*

*Proof.* See proof on [Lee, 2013, p. 10]. □

### 3.2.2 Smooth Manifolds

**Definition 124.** *A **coordinate chart** (or just a **chart**) on a topological manifold $M$ is a pair $(U, \varphi)$, where $U$ is an open subset of $M$ and $\varphi \colon U \to \widehat{U}$ is a homeomorphism from $U$ to an open subset $\widehat{U} = \varphi(U) \subseteq \mathbb{R}^n$ (see Figure 3.1.)*

**Definition 125.** *Let $M$ be a topological n-manifold. If $(U, \varphi), (V, \psi)$ are two charts such that $U \cap V \neq \varnothing$, the composite map $\psi \circ \varphi^{-1} \colon \varphi(U \cap V) \to \psi(U \cap V)$ is called the **transition map** from $\varphi$ to $\psi$ (see Figure 3.2 below). It is a composition of homeomorphisms, and is therefore itself a homeomorphism.*

Figure 3.1: A coordinate chart.

**Definition 126.** *An r-times continuously differentiable function $f: \mathbb{R}^n \to \mathbb{R}^m$ is called a $C^r$-***diffeomorphism*** if it is a bijection ($m = n$ and $f(x) = 0 \iff x = 0$) and its inverse $f^{-1}: \mathbb{R}^m \to \mathbb{R}^n$ is r-times continuously differentiable as well. For our purposes we will focus on the case $r = \infty$, in which case we call the $C^\infty$ function f a ***smooth function*** and instead of saying that it is a $C^\infty$ diffeomorphism we simply say that it is a ***diffeomorphism***. (These concepts will be extended to maps between more general manifolds in the next subsection.)*

**Definition 127.** *Two charts $(U, \varphi)$ and $(V, \psi)$ are said to be ***smoothly compatible*** if either $U \cap V = \varnothing$ or the transition map $\psi \circ \varphi^{-1}$ is a diffeomorphism. Since $\varphi(U \cap V)$ and $\psi(U \cap V)$ are open subsets of $\mathbb{R}^n$, smoothness of this map is to be interpreted in the ordinary sense of having continuous partial derivatives of all orders.*

**Definition 128.** *We define an ***atlas*** for a manifold M to be a collection of charts whose domains cover M. An atlas $\mathscr{A}$ is called a ***smooth atlas*** if any two charts in $\mathscr{A}$ are smoothly compatible with each other.*

Note that to show that an atlas is smooth, we need only verify that each transition map $\psi \circ \varphi^{-1}$ is smooth whenever $(U, \varphi)$ and $(V, \psi)$ are charts in $\mathscr{A}$; once we have proved this, it follows that $\psi \circ \varphi^{-1}$ is a diffeomorphism because its inverse $(\psi \circ \varphi^{-1})^{-1} = \varphi \circ \psi^{-1}$ is one of the transition maps we have already shown to be smooth. Alternatively, given two particular charts $(U, \varphi)$ and $(V, \psi)$ it is often easiest to show that they are smoothly compatible by verifying that $\psi \circ \varphi^{-1}$ is smooth and injective with nonsingular Jacobian at each point, and appealing to the following proposition:

Figure 3.2: A transition map.

**Proposition 38.** *Suppose $U \subseteq \mathbb{R}^n$ is an open subset, and $F\colon U \to \mathbb{R}^n$ is a smooth function whose Jacobian determinant is nonzero at every point in U. Then,*

  a) *F is an open map.*

  b) *if F is injective, then $F\colon U \to F(U)$ is a diffeomorphism.*

Our plan is to define a "smooth structure" on $M$ by giving a smooth atlas, and to define a function $f\colon M \to \mathbb{R}$ to be smooth if and only if $f \circ \varphi^{-1}$ is smooth in the sense of ordinary calculus for each coordinate chart $(U, \varphi)$ in the atlas. There is one minor technical problem with this approach: in general, there will be many possible atlases that give the "same" smooth structure, in that they all determine the same collection of smooth functions on $M$.

For example, consider the following pair of atlases on $\mathbb{R}^n$:

$$\mathcal{A}_\infty = \{(\mathbb{R}^n, \mathrm{Id}_{\mathbb{R}^n})\},$$
$$\mathcal{A}_\in = \{(B_1(x), \mathrm{Id}_{B_1(x)}) \mid x \in \mathbb{R}^n\}.$$

Although these are different smooth atlases, clearly a function $f \colon \mathbb{R}^n \to \mathbb{R}$ is smooth with respect to either atlas if and only if it is smooth in the sense of ordinary calculus. We could choose to define a smooth structure as an equivalence class of smooth atlases under an appropriate equivalence relation. However, it is more straightforward to make the following definitions:

**Definition 129.** *A smooth atlas $\mathscr{A}$ on M is **maximal** (or **complete**) if it is not properly contained in any larger smooth atlas. This just means that any chart that is smoothly compatible with every chart in $\mathscr{A}$ is already in $\mathscr{A}$.*

**Definition 130.** *If M is a topological manifold, a **smooth structure** on M is a maximal smooth atlas. A **smooth manifold** is a pair $(M, \mathscr{A})$, where M is a topological manifold and $\mathscr{A}$ is a smooth structure on M.*

It is generally not very convenient to define a smooth structure by explicitly describing a maximal smooth atlas, because such an atlas contains very many charts. Fortunately, we need only specify some smooth atlas, as the next proposition shows:

**Proposition 39.** *Let M be a topological manifold.*

  *a) Every smooth atlas $\mathscr{A}$ for M is contained in a unique maximal smooth atlas, called the smooth structure determined by $\mathscr{A}$.*

  *b) Two smooth atlases for M determine the same smooth structure if and only if their union is a smooth atlas.*

For example, if a topological manifold $M$ can be covered by a single chart, the smooth compatibility condition is trivially satisfied, so any such chart automatically determines a smooth structure on $M$. An example would be $\mathscr{A}_1 = \{(\mathbb{R}^n, \mathrm{Id}_{\mathbb{R}^n})\}$.

**Definition 131.** *A set $B \subseteq M$ is called a* **regular coordinate ball** *if there is a smooth coordinate ball $B' \supseteq \overline{B}$ and a smooth coordinate map $\varphi \colon B' \to \mathbb{R}^n$ such that for some positive real numbers $r < r'$, we have*

$$\varphi(B) = B_r(0), \quad \varphi(\overline{B}) = \overline{B}_r(0), \quad \text{and} \quad \varphi(B') = B_{r'}(0).$$

**Proposition 40.** *Every smooth manifold has a countable basis of regular coordinate balls.*

Usually we construct a smooth manifold structure in two stages: we start with a topological space and check that it is a topological manifold, and then we specify a smooth structure. However it is often more convenient to combine these two steps into a single construction, especially if we start with a set that is not already equipped with a topology. The following lemma provides a shortcut; it shows how, given a set with suitable "charts" that overlap smoothly, we can use the charts to define both a topology and a smooth structure on the set:

**Lemma 27** (**Smooth Manifold Chart Lemma**). *Let $M$ be a set, and suppose we are given a collection $\{U_\alpha\}$ of subsets of $M$ together with maps $\varphi_\alpha \colon U_\alpha \to \mathbb{R}^n$, such that the following properties are satisfied:*

a) *For each $\alpha$, $\varphi_\alpha$ is a bijection between $U_\alpha$ and an open subset $\varphi_\alpha(U_\alpha) \subseteq \mathbb{R}^n$.*

b) *For each $\alpha$ and $\beta$, the sets $\varphi_\alpha(U_\alpha \cap U_\beta)$ and $\varphi_\beta(U_\alpha \cap U_\beta)$ are open in $\mathbb{R}^n$.*

c) *Whenever $U_\alpha \cap U_\beta \neq \varnothing$, the map $\varphi_\beta \circ \varphi_\alpha^{-1} \colon \varphi_\alpha(U_\alpha \cap U_\beta) \to \varphi_\beta(U_\alpha \cap U_\beta)$ is smooth.*

d) *Countably many of the sets $U_\alpha$ cover $M$.*

e) *Whenever $p, q$ are distinct points in $M$, either there exists some $U_\alpha$ containing both $p$ and $q$ or there exist disjoint sets $U_\alpha, U_\beta$ with $p \in U_\alpha$ and $q \in U_\beta$.*

*Then $M$ has a unique smooth manifold structure such that each $(U_\alpha, \varphi_\alpha)$ is a smooth chart.*

*Proof.* See proof on [Lee, 2013, p. 22]. □

## 3.3 Smooth Maps

**Definition 132.** *Suppose M is a smooth n-manifold, k is a nonnegative integer, and $f: M \to \mathbb{R}^k$ is any function. We say that f is a **smooth function** if for every $p \in M$, there exists a smooth chart $(U, \varphi)$ for M whose domain contains p and such that the composite function $f \circ \varphi^{-1}$ is smooth on the open subset $\widehat{U} = \varphi(U) \subseteq \mathbb{R}^n$ (see Figure 3.3 below).*



Figure 3.3: Definition of smooth functions.

The most important special case is that of smooth real-valued functions $f: M \to \mathbb{R}$; the set of all such functions is denoted by $C^\infty(M)$. Because sums and constant multiples of smooth functions are smooth, it turns out that $C^\infty(M)$ is a vector space over $\mathbb{R}$.

The definition of smooth functions generalizes easily to maps between manifolds:

**Definition 133.** *Let M and N be smooth manifolds, and let $F: M \to N$ be any map. We say that F is a **smooth map** if for every $p \in M$, there exist smooth charts $(U, \varphi)$ containing p and $(V, \psi)$ containing $F(p)$ such that $F(U) \subseteq V$ and the composite map $\psi \circ F \circ \varphi^{-1}$ is smooth from $\varphi(U)$ to $\psi(V)$ (see Figure 3.4 below).*

Note that our previous definition of smoothness of real-valued or vector-valued functions can be viewed as a special case of this one, by taking $N = V = \mathbb{R}^k$ and $\psi = \mathrm{Id}: \mathbb{R}^k \to \mathbb{R}^k$.

Figure 3.4: Definition of smooth maps.

**Remark**: In spite of the apparent complexity of the definition, it is usually not hard to prove that a particular map is smooth. There are basically only three common ways to do so:

- Write the map in smooth local coordinates and recognize its component functions as compositions of smooth elementary functions.

- Exhibit the map as a composition of maps that are known to be smooth.

- Use some special-purpose theorem that applies to the particular case under consideration.

**Proposition 41.** *Every smooth map is continuous.*

*Proof.* Suppose $M$ and $N$ are smooth manifolds (with or without boundary,) and $F\colon M \to N$ is smooth. Given $p \in M$, smoothness of $F$ means there are smooth charts $(U, \varphi)$ containing $p$ and $(V, \psi)$ containing $F(p)$ such that $F(U) \subseteq V$ and $\psi \circ F \circ \varphi^{-1}\colon \varphi(U) \to \psi(V)$ is a smooth real/vector valued function, which is known to be continuous. Since $\varphi\colon U \to \varphi(U)$ and $\psi\colon V \to \psi(V)$ are homeomorphisms, this implies in turn that

$$F|_U = \psi^{-1} \circ (\psi \circ F \circ \varphi^{-1}) \circ \varphi\colon U \to V,$$

which is a composition of continuous maps. Since $F$ is continuous in a neighborhood of each point, it is continuous on $M$, as desired. □

**Proposition 42 (Equivalent Characterizations of Smoothness).** *Suppose M and N are smooth manifolds (with or without boundary,) and $F\colon M \to N$ is a map. Then F is smooth if and only if either of the following conditions is satisfied:*

a) *For every $p \in M$, there exist smooth charts $(U, \varphi)$ containing p and $(V, \psi)$ containing $F(p)$ such that $U \cap F^{-1}(V)$ is open in M and the composite map $\psi \circ F \circ \varphi^{-1}$ is smooth from $\varphi(U \cap F^{-1}(V))$ to $\psi(V)$.*

b) *F is continuous and there exist smooth atlases $\{(U_\alpha, \varphi_\alpha)\}$ and $\{(V_\beta, \psi_\beta)\}$ for M and N, respectively, such that for each $\alpha$ and $\beta$, $\psi_\beta \circ F \circ \varphi_\alpha^{-1}$ is a smooth map from $\varphi_\alpha(U_\alpha \cap F^{-1}(V_\beta))$ to $\psi_\beta(V_\beta)$.*

**Proposition 43 (Smoothness Is Local).** *Suppose M and N are smooth manifolds (with or without boundary), and let $F\colon M \to N$ be a map.*

a) *If every point $p \in M$ has a neighborhood U such that the restriction $F|_U$ is smooth, then F is smooth.*

b) *Conversely, if F is smooth, then its restriction to every open subset is smooth.*

The next corollary is essentially just a restatement of the previous proposition, but it gives a highly useful way of constructing smooth maps:

**Corollary 28 (Gluing Lemma for Smooth Maps).** *Let M and N be smooth manifolds (with or without boundary,) and let $\{U_\alpha\}_{\alpha \in A}$ be an open cover of M. Suppose that for each $\alpha \in A$, we are given a smooth map $F_\alpha\colon U_\alpha \to N$ such that the maps agree on overlaps, that is*

$$F_\alpha|_{U_\alpha \cap U_\beta} = F_\beta|_{U_\alpha \cap U_\beta} \quad \forall\, \alpha, \beta \in A.$$

*Then there exists a unique smooth map $F\colon M \to N$ such that $F|_{U_\alpha} = F_\alpha$ for each $\alpha \in A$.*

**Definition 134.** *If $F\colon M \to N$ is a smooth map, and $(U, \varphi)$ and $(V, \psi)$ are any smooth charts for M and N, respectively, we call the composite map $\widehat{F} = \psi \circ F \circ \varphi^{-1}$ the **coordinate representation** of F with respect to the given coordinates. It maps the set $\varphi(U \cap F^{-1}(V))$ to $\psi(V)$.*

**Proposition 44.** *Suppose $F: M \to N$ is a smooth map between smooth manifolds (with or without boundary). Then the coordinate representation of F with respect to every pair of smooth charts for M and N is smooth.*

**Remark**: As with real-valued or vector-valued functions, once we have chosen specific local coordinates in both the domain and codomain, we can often ignore the distinction between $\widehat{F}$ and $F$.

**Proposition 45.** *Let M, N, and P be smooth manifolds (with or without boundary.)*

a) *Every constant map $c: M \to N$ is smooth.*

b) *The identity map of M is smooth.*

c) *If $U \subseteq M$ is an open submanifold (with or without boundary), then the inclusion map $U \hookrightarrow M$ is smooth.*

d) *If $F: M \to N$ and $G: N \to P$ are smooth, then so is $G \circ F: M \to P$.*

**Proposition 46.** *Suppose $M_1, \ldots, M_k$ and N are smooth manifolds (with or without boundary), such that at most one of $M_1, \ldots, M_k$ has nonempty boundary. For each i, let $\pi_i: M_1 \times \cdots \times M_k \to M_i$ denote the projection onto the $M_i$ factor. Then a map $F: N \to M_1 \times \cdots \times M_k$ is smooth if and only if each of the component maps $F_i = \pi_i \circ F: N \to M_i$ is smooth.*

*Proof.* For simplicity we assume that $M_k$ is a smooth manifold with boundary. From Proposition 45 part d), it is clear that every $F_i$ is smooth if $F$ is smooth. Suppose that each $F_i$ is smooth and let $x \in N$. For each $i$, choose smooth charts $(U_i, \varphi_i)$ containing $x$ and $(V_i, \psi_i)$ containing $F_i(x)$ such that

$$\psi_i \circ F_i \circ \varphi_i^{-1} = \psi_i \circ \pi_i \circ F \circ \varphi_i^{-1}.$$

Now replace $U_1$ with $U = U_1 \cap \cdots \cap U_k$ and replace $\varphi_1$ with $\varphi_1|_U$; by Proposition 44, each map $\psi_i \circ \pi_i \circ F \circ \varphi_i^{-1}$ is smooth. Write $V = V_1 \times \cdots \times V_k$ and $\psi = \psi_1 \times \cdots \times \psi_k$ so that $(V, \psi)$ is a smooth chart containing $F(x)$. Since the $i^{th}$ component of $\psi \circ F \circ \varphi_1^{-1}$ is just $\psi_i \circ \pi_i \circ F \circ \varphi_1^{-1}$, the map $\psi \circ F \circ \varphi_1^{-1}$ is smooth. This shows that $F$ is smooth, as desired. $\square$

## 3.3.1   Partitions of Unity

The version of the gluing lemma for smooth maps that we presented on Corollary 28 is well defined on open subsets, but we cannot expect to glue together smooth maps defined on closed subsets and obtain a smooth result. For example, the two functions $f_+ \colon [0,\infty) \to \mathbb{R}$ and $f_- \colon (-\infty,0] \to \mathbb{R}$ defined by

$$f_+(x) = +x, \qquad x \in [0,\infty),$$
$$f_-(x) = -x, \qquad x \in (-\infty,0]$$

are both smooth and agree at the point $0$ where they overlap, but the continuous function $f \colon \mathbb{R} \to \mathbb{R}$ that they define, namely $f(x) = \|x\|$, is not smooth at the origin.

A disadvantage of this corollary is that in order to use it, we must construct maps that agree exactly on relatively large subsets of the manifold, which is too restrictive for some purposes. In this section we introduce partitions of unity, which are tools for "blending together" local smooth objects into global ones without necessarily assuming that they agree on overlaps.

All of our constructions in this section are based on the existence of smooth functions that are positive in a specified part of a manifold and identically zero in some other part. We begin by defining a smooth function on the real line that is zero for $t \leq 0$ and positive for $t > 0$:

**Lemma 28.** *The function $f \colon \mathbb{R} \to \mathbb{R}$ defined by*

$$f(t) = \begin{cases} e^{-1/t} & t > 0, \\ 0 & t \leq 0 \end{cases}$$

*is smooth (see* Figure 3.5 *below).*

**Lemma 29.** *Given any real numbers $r_1$ and $r_2$ such that $r_1 < r_2$, there exists a smooth function $h \colon \mathbb{R} \to \mathbb{R}$ such that*

$$h(t) = \begin{cases} 1 & t \leq r_1, \\ 0 < h(t) < 1 & r_1 < t < r_2, \\ 0 & t \geq r_2. \end{cases}$$

Figure 3.5: $f(t) = e^{-1/t}$.

*Proof.* Let $f$ be the function of the previous lemma, and set

$$h(t) = \frac{f(r_2 - t)}{f(r_2 - t) + f(t - r_1)}.$$



Figure 3.6: A cutoff function.

Such a function is usually called a ***cutoff function***. Note that the denominator is positive for all $t$, because at least one of the expressions $r_2 - t$ and $t - r_1$ is always positive. The desired properties of $h$ follow easily from those of $f$. $\qquad\square$

**Lemma 30.** *Given any positive real numbers $r_1 < r_2$, there is a smooth function $H\colon \mathbb{R}^n \to \mathbb{R}$ such*

*that*

$$H(x) = \begin{cases} 1 & x \in \overline{B}_{r_1}(0), \\ 0 < H(x) < 1 & x \in B_{r_2}(0) \smallsetminus \overline{B}_{r_1}(0), \\ 0 & x \in \mathbb{R}^n \smallsetminus B_{r_2}(0). \end{cases}$$

*Proof.* Just set $H(x) = h(\|x\|)$, where $h$ is the function of the preceding lemma. Clearly, $H$ is smooth on $\mathbb{R}^n \smallsetminus \{0\}$, because it is a composition of smooth functions there. Since it is identically equal to 1 on $B_{r_1}(0)$, it is smooth there too. $\qquad\square$

The function $H$ constructed in this lemma is an example of a ***smooth bump function***, a smooth real-valued function that is equal to 1 on a specified set and is zero outside a specified neighborhood of that set.

**Definition 135.** *Suppose $M$ is a topological space, and let $\chi = (X_\alpha)_{\alpha \in A}$ be an arbitrary open cover of $M$, indexed by a set $A$. A **partition of unity subordinate to** $\chi$ is an indexed family $(\psi_\alpha)_{\alpha \in A}$ of continuous functions $\psi_\alpha \colon M \to \mathbb{R}$ with the following properties:*

- $0 \leq \psi_\alpha \leq 1$ *for all $\alpha \in A$ and all $x \in M$.*

- $\operatorname{supp} \psi_\alpha \subseteq X_\alpha$ *for each $\alpha \in A$.*

- *The family of supports $(\operatorname{supp} \psi_\alpha)_{\alpha \in A}$ is locally finite.*

- $\sum_{\alpha \in A} \psi_\alpha(x) = 1$ *for all $x \in M$.*

*If $M$ is a smooth manifold with or without boundary (as opposed to just an arbitrary topological space,) a **smooth partition of unity** is one for which each of the functions $\psi_\alpha$ is smooth.*

**Theorem 103 (Existence of Partitions of Unity).** *Suppose $M$ is a smooth manifold (with or without boundary,) and $\chi = (X_\alpha)_{\alpha \in A}$ is any indexed open cover of $M$. Then there exists a smooth partition of unity subordinate to $\chi$.*

There are basically two different strategies for patching together locally defined smooth maps to obtain a global one. If you can define a map in a neighborhood of each point in such a way that the locally defined maps all agree where they overlap, then the local definitions

piece together to yield a global smooth map by *Corollary 28*. (This usually requires some sort of uniqueness result.) But if the local definitions are not guaranteed to agree, then you usually have to resort to a partition of unity. The trick then is showing that the patched-together objects still have the required properties.

**Definition 136.** *If M is a topological space, $A \subseteq M$ is a closed subset, and $U \subseteq M$ is an open subset containing A, a continuous function $\psi \colon M \to \mathbb{R}$ is called a **bump function for A supported in U** if $0 \leq \psi \leq 1$ on M, $\psi \equiv 1$ on A, and $\operatorname{supp} \psi \subseteq U$.*

**Proposition 47 (Existence of Smooth Bump Functions).** *Let M be a smooth manifold (with or without boundary). For any closed subset $A \subseteq M$ and any open subset U containing A, there exists a smooth bump function for A supported in U.*

**Definition 137.** *Suppose M and N are smooth manifolds (with or without boundary,) and $A \subseteq M$ is an arbitrary subset. If N has empty boundary, we say that a map $F \colon A \to N$ is **smooth on A** if it has a smooth extension in a neighborhood of each point: that is, if for every $p \in A$ there is an open subset $W \subseteq M$ containing p and a smooth map $\widetilde{F} \colon W \to N$ whose restriction to $W \cap A$ agrees with F. When $\partial N \neq \varnothing$, we say that $F \colon A \to N$ is **smooth on A** if for every $p \in A$ there exist an open subset $W \subseteq M$ containing p and a smooth chart $(V, \psi)$ for N whose domain contains $F(p)$, such that $F(W \cap A) \subseteq V$ and $\psi \circ F|_{W \cap A}$ is smooth as a map into $\mathbb{R}^n$ in the sense defined above (i.e., it has a smooth extension in a neighborhood of each point).*

**Lemma 31 (Extension Lemma for Smooth Functions).** *Suppose M is a smooth manifold (with or without boundary,) $A \subseteq M$ is a closed subset, and $f \colon A \to \mathbb{R}^k$ is a smooth function. For any open subset U containing A, there exists a smooth function $\widetilde{f} \colon M \to \mathbb{R}^k$ such that $\widetilde{f}|_A = f$ and $\operatorname{supp} \widetilde{f} \subseteq U$.*

**Remark**: The assumption in the extension lemma that the codomain of $f$ is $\mathbb{R}^k$, and not some other smooth manifold, is needed: for other codomains, extensions can fail to exist for topological reasons. For example, the identity map $\mathbb{S}^1 \to \mathbb{S}^1$ is smooth, but does not have even a continuous extension to a map from $\mathbb{R}^2$ to $\mathbb{S}^1$. Later on in the course we will show that a smooth map from a closed subset of a smooth manifold into a smooth manifold has a smooth extension if and only if it has a continuous one.

**Definition 138.** *If M is a topological space, an **exhaustion function for** $M$ is a continuous function $f: M \to \mathbb{R}$ with the property that the set $f^{-1}((-\infty, c])$ (called a **sublevel set of** $f$) is compact for each $c \in \mathbb{R}$.*

**Remark**: The name "exhaustion function" comes from the fact that as $n$ ranges over the positive integers, the sublevel sets $f^{-1}((-\infty, n])$ form an exhaustion of $M$ by compact sets; thus an exhaustion function provides a sort of continuous version of an exhaustion by compact sets. For example, the functions $f: \mathbb{R}^n \to \mathbb{R}$ and $g: \mathbb{B}^n \to \mathbb{R}$ given by

$$f(x) = \|x\|^2, \qquad g(x) = \frac{1}{1 - \|x\|^2}$$

are smooth exhaustion functions. Of course, if $M$ is compact, any continuous real-valued function on $M$ is an exhaustion function, so such functions are interesting only for noncompact manifolds.

**Proposition 48 (Existence of Smooth Exhaustion Functions).** *Every smooth manifold (with or without boundary) admits a smooth positive exhaustion function.*

The following theorem shows the remarkable fact that every closed subset of a manifold can be expressed as a level set of some smooth real-valued function:

**Theorem 104 (Level Sets of Smooth Functions).** *Let M be a smooth manifold. If C is any closed subset of M, there is a smooth nonnegative function $f: M \to \mathbb{R}$ such that $f^{-1}(0) = C$.*

## 3.4 Tangent Vectors

Before we get started, we should mention that we will be using the so called *Einstein Summation Convention* throughout all the following sections. It may be a bit weird at first, but once you become familiar with it the computations will look simpler (or at the very least

less cluttered). Because of the appearance of summations such as $\sum_i x^i E_i$ in this subject, we often abbreviate such a sum by omitting the summation sign:

$$\text{We write} \quad E(x) = x^i E_i \quad \text{as an abbreviation of} \quad E(x) = \sum_{i=1}^{n} x^i E_i.$$

The rule is the following: if the same index name (such as $i$ in the expression above) appears exactly twice in any monomial term, once as an upper index and once as a lower index, that term is understood to be summed over all possible values of that index, generally from 1 to the dimension of the space in question. The rule also applies in the following situation:

$$\text{We write} \quad E(x) = X^i \frac{\partial}{\partial x^i} \quad \text{as an abbreviation of} \quad E(x) = \sum_{i=1}^{n} X^i \frac{\partial}{\partial x^i}.$$

We point this out because note that in this latter case the index appears twice as a super-script; however since in $\partial/\partial x^i$ the index appears in the "denominator," it is counted as a subscript and our summation convention applies.

### 3.4.1  Geometric Tangent Vectors

**Definition 139.** *Given a fixed point $a \in \mathbb{R}^n$, let us define the **geometric tangent space to $\mathbb{R}^n$ at $a$**, denoted by $\mathbb{R}^n_a$, to be the set $\{a\} \times \mathbb{R}^n = \{(a, v) \mid v \in \mathbb{R}^n\}$. A **geometric tangent vector** in $\mathbb{R}^n$ is an element of $\mathbb{R}^n_a$ for some $a \in \mathbb{R}^n$.*

As a matter of notation, we abbreviate $(a, v)$ as $v_a$ (or sometimes $v|_a$ if it is clearer, for example if $v$ itself has a subscript). We think of $v_a$ as the vector $v$ with its initial point at $a$ (*Figure* 3.7).

**Remark:** Note that the geometric tangent space $\mathbb{R}^n_a$ is a real vector space under the natural operations

$$v_a + w_a = (v + w)_a \qquad \text{and} \qquad c(v_a) = (cv)_a \qquad \forall\, v_a, w_a \in \mathbb{R}^n_a, \quad \forall\, c \in \mathbb{R}.$$

One thing that a geometric tangent vector provides is a means of taking directional deriva-tives of functions. For example, any geometric tangent vector $v_a \in \mathbb{R}^n_a$ yields a map $D_v|_a \colon C^\infty(\mathbb{R}^n) \to \mathbb{R}$, which takes the directional derivative of a $C^\infty$ function $f$ in the direc-tion $v$ at $a$:

$$D_v\big|_a f = D_v f(a) = \frac{\mathrm{d}}{\mathrm{d}t}\Big|_{t=0} f(a + tv).$$

Figure 3.7: Geometric tangent space.

This operation is linear over $\mathbb{R}$ and satisfies the product rule:

$$D_v\big|_a(fg) = f(a)\, D_v\big|_a g + g(a)\, D_v\big|_a f.$$

If $v_a = v^i e_i\big|_a$ (note the Einstein summation convention that we previously discussed) in terms of the standard basis, then by the chain rule $D_v\big|_a f$ can be written more concretely as

$$D_v\big|_a f = v^i \frac{\partial f}{\partial x^i}(a).$$

With this construction in mind, we make the following definition:

**Definition 140.** *If $a$ is a point of $\mathbb{R}^n$, a map $\omega \colon C^\infty(\mathbb{R}^n) \to \mathbb{R}$ is called a **derivation at** $a$ if it is linear over $\mathbb{R}$ and satisfies the product rule $\omega(fg)(a) = f(a)\omega(g) + g(a)\omega(f)$.*

**Remark:** Note that the set of all derivations of $C^\infty(\mathbb{R}^n)$ at $a$, denoted $\mathfrak{D}_a\mathbb{R}^n$, is a real vector space under the operations

$$(\omega_1 + \omega_2)(f) = \omega_1(f) + \omega_2(f) \qquad \text{and} \qquad (c\,\omega)(f) = c(\omega(f)).$$

The most important (and perhaps somewhat surprising) fact about $\mathfrak{D}_a\mathbb{R}^n$ is that it is finite-dimensional, and in fact is naturally isomorphic to the geometric tangent space $\mathbb{R}^n_a$ that we defined above. The proof is based on the following lemma:

**Lemma 32 (Properties of Derivations).** *Suppose $a \in \mathbb{R}^n$, $\omega \in \mathfrak{D}_a\mathbb{R}^n$, and $f, g \in C^\infty(\mathbb{R}^n)$. Then we have the following:*

*a)* If $f$ is a constant function, then $\omega(f) = 0$.

*b)* If $f(a) = g(a) = 0$, then $\omega(fg) = 0$.

The next proposition shows that derivations at $a$ are in one-to-one correspondence with geometric tangent vectors:

**Proposition 49.** *Let $a \in \mathbb{R}^n$. Then,*

*a)* *For each geometric tangent vector $v_a \in \mathbb{R}_a^n$, the map $D_v\big|_a \colon C^\infty(\mathbb{R}^n) \to \mathbb{R}$ is a derivation at a.*

*b)* *The map $v_a \mapsto D_v\big|_a$ is an isomorphism from $\mathbb{R}_a^n$ onto $\mathfrak{D}_a\mathbb{R}^n$.*

**Corollary 29.** *For any $a \in \mathbb{R}^n$, the n derivations*

$$\frac{\partial}{\partial x^1}\bigg|_a, \cdots, \frac{\partial}{\partial x^n}\bigg|_a \quad \text{defined by} \quad \frac{\partial}{\partial x^i}\bigg|_a f = \frac{\partial f}{\partial x^i}(a)$$

*form a basis for $\mathfrak{D}_a\mathbb{R}^n$, which therefore has dimension n.*

*Proof.* Apply the previous proposition and note that $\partial/\partial x^i\big|_a = D_{e_i}\big|_a$. $\qquad\square$

### 3.4.2 Tangent Vectors on Manifolds

Now we are in a position to define tangent vectors on manifolds:

**Definition 141.** *Let M be a smooth manifold (with or without boundary), and let p be a point of M. A linear map. A linear map $v \colon C^\infty(M) \to \mathbb{R}$ is called a **derivation at** p if it satisfies the product rule*

$$v(fg)(p) = f(p)v(g) + g(p)v(f) \quad \text{for all } f, g \in C^\infty(M).$$

*The set of all derivations of $C^\infty(M)$ at p, denoted by $T_pM$, is a vector space called the **tangent space to** M **at** p. An element of $T_pM$ is called a **tangent vector at** p.*

The following lemma is the analogue of Lemma 32 for manifolds:

**Lemma 33 (Properties of Tangent Vectors on Manifolds).** *Suppose M is a smooth manifold (with or without boundary), $p \in M$, $v \in T_pM$, and $f, g \in C^\infty(M)$. Then we have the following:*

   *a) If $f$ is a constant function, then $v(f) = 0$.*

   *b) If $f(p) = g(p) = 0$, then $v(fg) = 0$.*

To relate the abstract tangent spaces we have defined on manifolds to geometric tangent spaces in $\mathbb{R}^n$, we have to explore the way smooth maps affect tangent vectors. In the case of a smooth map between Euclidean spaces, the total derivative of the map at a point (represented by its Jacobian matrix) is a linear map that represents the "best linear approximation" to the map near the given point. In the manifold case there is a similar linear map, but it makes no sense to talk about a linear map between manifolds. Instead, it will be a linear map between tangent spaces.

**Definition 142.** *If M and N are smooth manifolds (with or without boundary) and $F\colon M \to N$ is a smooth map, then for each $p \in M$ we define a map*

$$\mathrm{d}F_p\colon T_pM \to T_{F(p)}N,$$

*called the **differential of F at** $p$, as follows: Given $v \in T_pM$, we let $\mathrm{d}F_p(v)$ be the derivation at $F(p)$ that acts on $f \in C^\infty(N)$ by the rule $\mathrm{d}F_p(v)(f) = v(f \circ F)$ (see Figure 3.8 below).*



Figure 3.8: The differential.

Note that if $f \in C^\infty(N)$, then $f \circ F \in C^\infty(M)$ so $v(f \circ F)$ makes sense. The operator $\mathrm{d}F_p(v)\colon C^\infty(N) \to \mathbb{R}$ is linear because $v$ is, and is a derivation at $F(p)$ because for any

$f, g \in C^\infty(N)$, we have

$$
\begin{aligned}
dF_p(v)(fg) &= v\left((fg) \circ F\right) = v\left((f \circ F)(g \circ F)\right) \\
&= f \circ F(p)\, v(g \circ F) + g \circ F(p)\, v(f \circ F) \\
&= f\left(F(p)\right) dF_p(v)(g) + g\left(F(p)\right) dF_p(v)(f).
\end{aligned}
$$

**Proposition 50 (Properties of Differentials).** *Let M, N, and S be smooth manifolds (with or without boundary), let $F\colon M \to N$ and $G\colon N \to S$ be smooth maps, and let $p \in M$.*

  a)  $dF_p\colon T_pM \to T_{F(p)}N$ *is linear.*

  b)  $d(G \circ F)_p = dG_{F(p)} \circ dF_p\colon T_pM \to T_{G\circ F(p)}S.$

  c)  $d(\mathrm{Id}_M)_p = \mathrm{Id}_{T_p(M)}\colon T_pM \to T_pM.$

  d)  *If F is a diffeomorphism, then $dF_p\colon T_pM \to T_{F(p)}N$ is an isomorphism, and $(dF_p)^{-1} = d(F^{-1})_{F(p)}.$*

The next proposition indicates that tangent vectors act locally.

**Proposition 51.** *Let M be a smooth manifold (with or without boundary), $p \in M$, and $v \in T_pM$. If $f, g \in C^\infty(M)$ agree on some neighborhood of p, then $v(f) = v(g)$.*

Using this proposition, we can identify the tangent space to an open submanifold with the tangent space to the whole manifold:

**Proposition 52 (The Tangent Space to an Open Submanifold).** *Let M be a smooth manifold (with or without boundary), let $U \subseteq M$ be an open subset, and let $\iota\colon U \hookrightarrow M$ be the inclusion map. For every $p \in U$, the differential $d\iota_p\colon T_pU \to T_pM$ is an isomorphism.*

**Remark:** Given an open subset $U \subseteq M$, the isomorphism $d\iota_p$ between $T_pU$ and $T_pM$ is canonically defined, independently of any choices. Hence from now on we identify $T_pU$ with $T_pM$ for any point $p \in U$.

**Proposition 53 (Dimension of the Tangent Space).** *If M is an n-dimensional smooth manifold, then for each $p \in M$, the tangent space $T_pM$ is an n-dimensional vector space.*

Recall that every finite-dimensional vector space has a natural smooth manifold structure that is independent of any choice of basis or norm. The following proposition shows that the tangent space to a vector space can be naturally identified with the vector space itself. Suppose $V$ is a finite-dimensional vector space and $a \in V$. Just as we did earlier in the case of $\mathbb{R}^n$, for any vector $v \in V$, we define a map $D_v\big|_a : C^\infty(V) \to \mathbb{R}$ by

$$D_v\big|_a f = \frac{d}{dt}\bigg|_{t=0} f(a + tv). \tag{3.2}$$

**Proposition 54 (The Tangent Space to a Vector Space).** *Suppose V and W are finite-dimensional vector spaces with their respective standard smooth manifold structures. For each point $a \in V$, the map $v \mapsto D_v\big|_a$ defined by (3.2) is a canonical isomorphism from V to $T_aV$, such that for any linear map $L : V \to W$, the following diagram commutes:*

$$
\begin{array}{ccc}
V & \xrightarrow{\;\cong\;} & T_aV \\
{\scriptstyle L}\downarrow & & \downarrow{\scriptstyle dL_a} \\
W & \xrightarrow[\;\cong\;]{} & T_{L_a}W
\end{array}
$$

It is important to understand that each isomorphism $V \cong T_aV$ is canonically defined, independently of any choice of basis. Because of this result, we can routinely identify tangent vectors to a finite-dimensional vector space with elements of the space itself. More generally, if $M$ is an open submanifold of a vector space $V$, we can combine our identifications $T_pM \leftrightarrow T_pV \leftrightarrow V$ to obtain a canonical identification of each tangent space to $M$ with $V$. For example, since $GL(n, \mathbb{R})$ is an open submanifold of the vector space $M(n, \mathbb{R})$, we can identify its tangent space at each point (i.e., matrix) $X \in GL(n, \mathbb{R})$ with the full space of matrices $M(n, \mathbb{R})$.

There is another natural identification for tangent spaces to a product manifold:

**Proposition 55 (The Tangent Space to a Product Manifold).** *Let $M_1, \ldots, M_k$ be smooth manifolds, and for each $j$, let $\pi_j \colon M_1 \times \cdots \times M_k \to M_j$ be the projection onto the $M_j$ factor. For any point $p = (p_1, \ldots, p_k) \in M_1 \times \cdots \times M_k$ and tangent vector $v \in T_p(M_1 \times \cdots \times M_k)$, the map*

$$\alpha \colon T_p(M_1 \times \cdots \times M_k) \longrightarrow T_{p_1}M_1 \oplus \cdots \oplus T_{p_k}M_k$$

*defined by*

$$\alpha(v) = \big(\mathrm{d}(\pi_1)_p(v), \ldots, \mathrm{d}(\pi_k)_p(v)\big) \tag{3.3}$$

*is an isomorphism. The same is true if one of the spaces $M_i$ is a smooth manifold with boundary.*

Once again, because the isomorphism (3.3) is canonically defined, independently of any choice of coordinates, we can consider it as a canonical identification, and we will always do so. Thus, for example, we identify $T_{(p,q)}(M \times N)$ with $T_pM \oplus T_qN$, and treat $T_pM$ and $T_qN$ as subspaces of $T_{(p,q)}(M \times N)$.

### 3.4.3   Computations in Coordinates

Suppose $M$ is a smooth manifold and let $(U, \varphi)$ be a smooth coordinate chart on $M$. Then $\varphi$ is, in particular, a diffeomorphism from $U$ to an open subset $\widehat{U} \subseteq \mathbb{R}^n$. Combining Propositions 52 and 50 part d) from above, we see that $d\varphi(p) \colon T_pM \to \mathfrak{D}_{\varphi(p)}\mathbb{R}^n$ is an isomorphism. Then by Corollary 29, the derivations $\partial/\partial x^1|_{\varphi(p)}, \ldots, \partial/\partial x^n|_{\varphi(p)}$ form a basis for $\mathfrak{D}_{\varphi(p)}\mathbb{R}^n$. Therefore, the preimages of these vectors under the isomorphism $d\varphi_p$ form a basis for $T_pM$.

In keeping with our standard practice of treating coordinate maps as identifications whenever possible, we use the notation $\partial/\partial x^i|_p$ for these vectors, characterized by either of the following expressions:

$$\frac{\partial}{\partial x^i}\bigg|_p = (\mathrm{d}\varphi_p)^{-1}\left(\frac{\partial}{\partial x^i}\bigg|_{\varphi(p)}\right) = \mathrm{d}(\varphi^{-1})_{\varphi(p)}\left(\frac{\partial}{\partial x^i}\bigg|_{\varphi(p)}\right).$$

Unwinding the definitions, we see that $\partial/\partial x^i|_p$ acts on a function $f \in C^\infty(U)$ by

$$\frac{\partial}{\partial x^i}\bigg|_p f = \frac{\partial}{\partial x^i}\bigg|_{\varphi(p)}(f \circ \varphi^{-1}) = \frac{\partial \widehat{f}}{\partial x^i}(\widehat{p}),$$

where $\hat{f} = f \circ \varphi^{-1}$ is the coordinate representation of $f$, and $\hat{p} = (p^1, \ldots, p^n) = \varphi(p)$ is the coordinate representation of $p$. In other words, $\partial/\partial x^i|_p$ is just the derivation that takes the $i^{th}$ partial derivative of (the coordinate representation of) $f$ at (the coordinate representation of) $p$. The vectors $\partial/\partial x^i|_p$ are called the ***coordinate vectors at $p$*** associated with the given coordinate system. In the special case of standard coordinates on $\mathbb{R}^n$, the vectors $\partial/\partial x^i|_p$ are literally just the partial derivative operators.

The following proposition summarizes the discussion so far:

**Proposition 56.** *Let $M$ be a smooth n-manifold (with or without boundary), and let $p \in M$. Then $T_pM$ is an n-dimensional vector space, and for any smooth chart $(U, (x^i))$ containing $p$, the coordinate vectors $\partial/\partial x^1|_p, \ldots, \partial/\partial x^n|_p$ form a basis for $T_pM$.*

Thus, a tangent vector $v \in T_pM$ can be written uniquely as a linear combination

$$v = v^i \frac{\partial}{\partial x^i}\Big|_p,$$

The ordered basis $\left(\partial/\partial x^i|_p\right)$ is called a ***coordinate basis for $T_pM$***, and the coefficients $(v^1, \ldots, v^n)$ are called the ***components of $v$*** with respect to the coordinate basis. If $v$ is known, its components can be computed easily from its action on the coordinate functions. For each $j$, the components of $v$ are given by $v^j = v(x^j)$ (where we think of $x^j$ as a smooth real-valued function on $U$), because

$$v(x^j) = \left( v^i \frac{\partial}{\partial x^i}\Big|_p \right)(x^j) = v^i \frac{\partial x^j}{\partial x^i}(p) = v^j.$$

### 3.4.4 The Differential in Coordinates

Next we explore how differentials look in coordinates. We begin by considering the special case of a smooth map $F \colon U \to V$, where $U \subseteq \mathbb{R}^n$ and $V \subseteq \mathbb{R}^m$ are open subsets of Euclidean spaces. For any $p \in U$, we will determine the matrix of $dF_p \colon T_p\mathbb{R}^n \to T_{F(p)}\mathbb{R}^m$ in terms of the standard coordinate bases. Using $(x^1, \ldots, x^n)$ to denote the coordinates in the domain and $(y^1, \ldots, y^m)$ to denote those in the codomain, and letting $f \in C^\infty(\mathbb{R}^m)$, we use the chain

rule to compute the action of $dF_p$ on a typical basis vector as follows:

$$dF_p \left( \frac{\partial}{\partial x^i} \bigg|_p \right) f = \frac{\partial}{\partial x^i} \bigg|_p (f \circ F)$$

$$= \frac{\partial f}{\partial y^j}(F(p)) \frac{\partial F^j}{\partial x^i}(p)$$

$$= \left( \frac{\partial F^j}{\partial x^i}(p) \frac{\partial}{\partial y^j} \bigg|_{F(p)} \right) f.$$

Thus we have that

$$dF_p \left( \frac{\partial}{\partial x^i} \bigg|_p \right) = \frac{\partial F^j}{\partial x^i}(p) \frac{\partial}{\partial y^j} \bigg|_{F(p)}.$$

In other words, the matrix of $dF_p$ in terms of the coordinate bases is

$$\begin{pmatrix} \frac{\partial F^1}{\partial x^1}(p) & \cdots & \frac{\partial F^1}{\partial x^n}(p) \\ \vdots & \ddots & \vdots \\ \frac{\partial F^m}{\partial x^1}(p) & \cdots & \frac{\partial F^m}{\partial x^n}(p) \end{pmatrix}.$$

Recall that the columns of the matrix are the components of the images of the basis vectors. This matrix is none other than the Jacobian matrix of $F$ at $p$, which is the matrix representation of the total derivative $DF(p) \colon \mathbb{R}^n \to \mathbb{R}^m$. Therefore, in this case, $dF_p \colon T_p\mathbb{R}^n \to T_{F(p)}\mathbb{R}^m$ corresponds to the total derivative under our usual identification of Euclidean spaces with their tangent spaces.

Now consider the more general case of a smooth map $F \colon M \to N$ between smooth manifolds (with or without boundary). Choosing smooth coordinate charts $(U, \varphi)$ for $M$ containing $p$ and $(V, \psi)$ for $N$ containing $F(p)$, we obtain the coordinate representation

$$\widehat{F} = \psi \circ F \circ \varphi^{-1} \colon \varphi(U \cap F^{-1}(V)) \to \psi(V) \qquad \text{(see Figure 3.9).}$$

Now let $\widehat{p} = \varphi(p)$ denote the coordinate representation of $p$. By the above computation, $d\widehat{F}_{\widehat{p}}$ is represented with respect to the standard coordinate bases by the Jacobian matrix of

Figure 3.9: The differential in coordinates.

$\widehat{F}$ at $\widehat{p}$. Using the fact that $F \circ \varphi^{-1} = \psi^{-1} \circ \widehat{F}$, we compute

$$
\begin{aligned}
\mathrm{d}F_p\left(\left.\frac{\partial}{\partial x^i}\right|_p\right) &= \mathrm{d}F_p\left(\mathrm{d}(\varphi^{-1})_{\widehat{p}}\left(\left.\frac{\partial}{\partial x^i}\right|_{\widehat{p}}\right)\right) \\
&= \mathrm{d}(\psi^{-1})_{\widehat{F}(\widehat{p})}\left(\mathrm{d}\widehat{F}_{\widehat{p}}\left(\left.\frac{\partial}{\partial x^i}\right|_{\widehat{p}}\right)\right) \\
&= \mathrm{d}(\psi^{-1})_{\widehat{F}(\widehat{p})}\left(\frac{\partial \widehat{F}^j}{\partial x^i}(\widehat{p})\left.\frac{\partial}{\partial y^j}\right|_{\widehat{F}(\widehat{p})}\right) \\
&= \frac{\partial \widehat{F}^j}{\partial x^i}(\widehat{p})\left.\frac{\partial}{\partial y^j}\right|_{F(p)}.
\end{aligned}
$$

Thus, $\mathrm{d}F_p$ is represented in coordinate bases by the Jacobian matrix of (the coordinate representative of) $F$. In fact, the definition of the differential was cooked up precisely to give a coordinate-independent meaning to the Jacobian matrix. In the differential geometry literature, the differential is sometimes called the *tangent map*, the *total derivative*, or simply the *derivative* of $F$. Because it "pushes" tangent vectors forward from the domain manifold to the codomain, it is also called the *(pointwise) pushforward*.

### 3.4.5 The Tangent Bundle

**Definition 143.** *Given a smooth manifold M (with or without boundary), we define the **tangent bundle of M**, denoted by TM, to be the disjoint union of the tangent spaces at all points of M:*

$$TM = \amalg_{p \in M} T_p M.$$

*We usually write an element of this disjoint union as an ordered pair $(p, v)$, with $p \in M$ and $v \in T_p M$. The tangent bundle comes equipped with a natural **projection map** $\pi \colon TM \to M$, which sends each vector v in $T_p M$ to the point p at which it is tangent: $\pi(p, v) = p$.*

**Proposition 57.** *For any smooth n-manifold M, the tangent bundle TM has a natural topology and smooth structure that make it into a 2n-dimensional smooth manifold. With respect to this structure, the projection $\pi \colon TM \to M$ is smooth.*

*Proof.* You can find the proof on [Lee, 2013, p. 66]. □

**Proposition 58.** *If M is a smooth n-manifold (with or without boundary), and M can be covered by a single smooth chart, then TM is diffeomorphic to $M \times \mathbb{R}^n$.*

**Definition 144.** *By putting together the differentials of F at all points of M, we obtain a globally defined map between tangent bundles, called the **global differential** or **global tangent map** and denoted by $dF \colon TM \to TN$. This is just the map whose restriction to each tangent space $T_p M \subseteq TM$ is $dF_p$ (when we apply the differential of F to a specific vector $v \in T_p M$, we can write either $dF_p(v)$ or $dF(v)$, depending on how much emphasis we wish to give to the point p).*

**Proposition 59.** *If $F \colon M \to N$ is a smooth map, then its global differential $dF \colon TM \to TN$ is a smooth map.*

The following properties of the global differential follow immediately from Proposition 50:

**Corollary 30 (Properties of the Global Differential).** *Let M, N, and S be smooth manifolds (with or without boundary), let $F \colon M \to N$ and $G \colon N \to S$ be smooth maps, and let $p \in M$. Then,*

a) $\mathrm{d}(G \circ F) = \mathrm{d}G \circ \mathrm{d}F$.

b) $\mathrm{d}(\mathrm{Id}_M) = \mathrm{Id}_{TM}$.

c) *If $F$ is a diffeomorphism, then $\mathrm{d}F\colon TM \to TN$ is also a diffeomorphism, and $(\mathrm{d}F)^{-1} = \mathrm{d}(F^{-1})$.*

## 3.5 Submersions, Immersions, Embeddings

In this section we discuss a corollary of the *Inverse Function Theorem*, known as the *Rank Theorem*, and some of its most crucial ramifications. We then delve more deeply into smooth embeddings and smooth submersions, and apply the theory to a particularly useful class of smooth submersions: the smooth covering maps.

### 3.5.1 Maps of Constant Rank

**Definition 145.** *Suppose $M$ and $N$ are smooth manifolds (with or without boundary). Given a smooth map $F\colon M \to N$ and a point $p \in M$, we define the **rank of $F$ at $p$** to be the rank of the linear map $\mathrm{d}F_p\colon T_pM \to T_{F(p)}N$; it is the rank of the Jacobian matrix of $F$ in any smooth chart, or the dimension of $\mathrm{Im}(\mathrm{d}F_p) \subseteq T_{F(p)}N$. If $F$ has the same rank $r$ at every point, we say that it has **constant rank**, and write $\mathrm{rank}\, F = r$.*

Because the rank of a linear map is never higher than the dimension of either its domain or its codomain, the rank of $F$ at each point is bounded above by $\min\{\dim M, \dim N\}$. If the rank of $\mathrm{d}F_p$ is equal to this upper bound, we say that *$F$ **has full rank at** $p$*, and if $F$ has full rank everywhere, we say *$F$ **has full rank**.*

The most important constant-rank maps are those of full rank:

**Definition 146.** *A smooth map $F\colon M \to N$ is called a **smooth submersion** if its differential is surjective at each point (or equivalently, if $\operatorname{rank} F = \dim N$). It is called a **smooth immersion** if its differential is injective at each point (equivalently, $\operatorname{rank} F = \dim M$).*

**Proposition 60.** *Suppose $F\colon M \to N$ is a smooth map and $p \in M$. If $dF_p$ is surjective, then $p$ has a neighborhood $U$ such that $F|_U$ is a submersion. If $dF_p$ is injective, then $p$ has a neighborhood $U$ such that $F|_U$ is an immersion.*

**Example 54.** *a) Suppose $M_1, \ldots, M_k$ are smooth manifolds. Then each of the projection maps $\pi_i\colon M_1 \times \cdots \times M_k \to M_i$ is a smooth submersion. In particular, the projection $\pi\colon \mathbb{R}^{n+k} \to \mathbb{R}^n$ onto the first $n$ coordinates is a smooth submersion.*

*b) If $\gamma\colon I \to M$ is a smooth curve in a smooth manifold $M$ (with or without boundary), then $\gamma$ is a smooth immersion if and only if $\gamma'(t) \neq 0$ for all $t \in I$.*

*c) If $M$ is a smooth manifold and its tangent bundle $TM$ is given the smooth manifold structure described in the proof to Proposition 57 (see [**?**, 66]), the projection $\pi\colon TM \to M$ is a smooth submersion. To verify this, just note that with respect to any smooth local coordinates $(x^i)$ on an open subset $U \subseteq M$ and the corresponding natural coordinates $(x^i, v^i)$ on $\pi^{-1}(U) \subseteq TM$, the coordinate representation of $\pi$ is $\widehat{\pi}(x, v) = x$.*

*d) The smooth map $X\colon \mathbb{R}^2 \to \mathbb{R}^3$ given by*

$$X(u, v) = ((2 + \cos 2\pi u) \cos 2\pi v, (2 + \cos 2\pi u) \sin 2\pi v, \sin 2\pi u)$$

*is a smooth immersion of $\mathbb{R}^2$ into $\mathbb{R}^3$ whose image is the doughnut-shaped surface obtained by revolving the circle $(y - 2)^2 + z^2 = 1$ in the $(y, z)$-plane about the z-axis (see Figure 3.10 below).* ✿

**Definition 147.** *If $M$ and $N$ are smooth manifolds (with or without boundary), a map $F\colon M \to N$ is called a **local diffeomorphism** if every point $p \in M$ has a neighborhood $U$ such that $F(U)$ is open in $N$ and $F|_U\colon U \to F(U)$ is a diffeomorphism.*

The next theorem is the key to the most important properties of local diffeomorphisms:

Figure 3.10: A torus of revolution in $\mathbb{R}^3$.

**Theorem 105 (Inverse Function Theorem for Manifolds).** *Suppose M and N are smooth manifolds (without boundary)[1] and $F\colon M \to N$ is a smooth map. If $p \in M$ is a point such that $\mathrm{d}F_p$ is invertible, then there are connected neighborhoods $U_0$ of $p$ and $V_0$ of $F(p)$ such that $F|_{U_0}\colon U_0 \to V_0$ is a diffeomorphism.*

**Proposition 61 (Elementary Properties of Local Diffeomorphisms).** *We have the following properties for local diffeomorphisms:*

   *a)* *Every composition of local diffeomorphisms is a local diffeomorphism.*

   *b)* *Every finite product of local diffeomorphisms between smooth manifolds is a local diffeomorphism.*

   *c)* *Every local diffeomorphism is a local homeomorphism and an open map.*

   *d)* *The restriction of a local diffeomorphism to an open submanifold (with or with out boundary) is a local diffeomorphism.*

   *e)* *Every diffeomorphism is a local diffeomorphism.*

   *f)* *Every bijective local diffeomorphism is a diffeomorphism.*

   *g)* *A map between smooth manifolds (with or without boundary) is a local diffeomorphism if and only if in a neighborhood of each point of its domain, it has a coordinate representation that is a local diffeomorphism.*

---

[1] To see a case where the theorem fails for a map whose domain has nonempty boundary, see Problem $4-1$ on [Lee, 2013, p. 95].

**Proposition 62.** *Suppose $M$ and $N$ are smooth manifolds (without boundary), and $F\colon M \to N$ is a map. Then we have the following:*

a) *$F$ is a local diffeomorphism if and only if it is both a smooth immersion and a smooth submersion.*

b) *If $\dim M = \dim N$ and $F$ is either a smooth immersion or a smooth submersion, then it is a local diffeomorphism.*

The most important fact about constant-rank maps is the following consequence of the inverse function theorem, which says that a constant-rank smooth map can be placed locally into a particularly simple canonical form by a change of coordinates. It is a nonlinear version of the canonical form theorem for linear maps (see [Lee, 2013, p. 626, Theorem $B$.20]):

**Theorem 106 (Rank Theorem).** *Suppose $M$ and $N$ are smooth manifolds of dimensions $m$ and $n$, respectively, and $F\colon M \to N$ is a smooth map with constant rank $r$. For each $p \in M$, there exist smooth charts $(U, \varphi)$ for $M$ centered at $p$ and $(V, \psi)$ for $N$ centered at $F(p)$ such that $F(U) \subseteq V$, in which $F$ has a coordinate representation of the form*

$$\widehat{F}(x^1, \ldots, x^r, x^{r+1}, \ldots, x^m) = (x^1, \ldots, x^r, 0, \ldots, 0).$$

*In particular, if $F$ is a smooth submersion, this becomes*

$$\widehat{F}(x^1, \ldots, x^n, x^{n+1}, \ldots, x^m) = (x^1, \ldots, x^n),$$

*and if $F$ is a smooth immersion, it is*

$$\widehat{F}(x^1, \ldots, x^m) = (x^1, \ldots, x^m, 0, \ldots, 0).$$

The next corollary can be viewed as a more invariant statement of the rank theorem. It says that constant-rank maps are precisely the ones whose local behavior is the same as that of their differentials:

**Corollary 31.** *Let $M$ and $N$ be smooth manifolds, let $F\colon M \to N$ be a smooth map, and suppose $M$ is connected. Then the following are equivalent:*

a) *For each $p \in M$ there exist smooth charts containing $p$ and $F(p)$, in which the coordinate representation of $F$ is linear.*

b) *F has constant rank.*

The rank theorem is a purely local statement. However, it has the following powerful global consequence.

**Theorem 107 (Global Rank Theorem).** *Let M and N be smooth manifolds and suppose $F: M \to N$ is a smooth map of constant rank.*

a) *If F is surjective, it is a smooth submersion.*

b) *If F is injective, it is a smooth immersion.*

c) *If F is bijective, it is a diffeomorphism.*

## 3.5.2  Embeddings

One special kind of immersion is particularly important:

**Definition 148.** *If M and N are smooth manifolds (with or without boundary), a **smooth embedding of M into N** is a smooth immersion $F: M \to N$ that is also a topological embedding, i.e., a homeomorphism onto its image $F(M) \subseteq N$ in the subspace topology.*

Remark: Note that a smooth embedding is a map that is both a topological embedding and a smooth immersion, not just a topological embedding that happens to be smooth.

**Example 55 (Smooth Embeddings).** a) *If M is a smooth manifold (with or without boundary) and $U \subseteq M$ is an open submanifold, the inclusion map $U \hookrightarrow M$ is a smooth embedding.*

b) *If $M_1, \dots, M_k$ are smooth manifolds and $p_i \in M_i$ are arbitrarily chosen points, each of the maps $\iota_j: M_j \to M_1 \times \cdots \times M_k$ given by*

$$\iota_j(q) = (p_1, \dots, p_{j-1}, q, p_{j+1}, \dots, p_k)$$

*is a smooth embedding. In particular, the inclusion map $\mathbb{R}^n \hookrightarrow \mathbb{R}^{n+k}$ given by*

$$(x^1, \dots, x^n) \hookrightarrow (x^1, \dots, x^n, \underbrace{0, \dots, 0}_{k-n \ zeroes})$$

*is a smooth embedding.*

To understand more fully what it means for a map to be a smooth embedding, it is useful to bear in mind some examples of injective smooth maps that are <u>NOT</u> smooth embeddings. The next three examples illustrate three rather different ways in which this can happen:

**Example 56 (Smooth Topological Embedding).** *The map $\gamma \colon \mathbb{R} \to \mathbb{R}^2$ given by $\gamma(t) = (t^3, 0)$ is a smooth map and a topological embedding. However, since $\gamma'(0) = 0$, we have that $\gamma$ is not a smooth immersion and thus not a smooth embedding.* ☕

**Example 57 (The Figure-Eight Curve).** *Consider the curve $\beta \colon (-\pi, \pi) \to \mathbb{R}^2$ defined by*

$$\beta(t) = (\sin 2t, \sin t).$$

*Its image is a set that looks like a figure-eight in the plane (see* Figure 3.11*), sometimes called a* **lemniscate**. *(It is the locus of points $(x, y)$ where $x^2 = 4y^2(1 - y^2)$, as you can check.) It is easy to*



Figure 3.11: A figure-eight curve (lemniscate).

*see that $\beta$ is an injective smooth immersion because $\beta'(t)$ never vanishes; but it is not a topological embedding, because its image is compact in the subspace topology, while its domain (the open set $(-\pi, \pi)$) is not.* ☕

**Example 58 (A Dense Curve on the Torus).** *Let $\mathbb{T}^2 = \mathbb{S}^1 \times \mathbb{S}^1 \subseteq \mathbb{C}^2$ denote the torus, and let $\alpha$ be any irrational number. The map $\gamma \colon \mathbb{R} \to \mathbb{T}^2$ given by*

$$\gamma(t) = \left( e^{2\pi i t}, e^{2\pi i \alpha t} \right)$$

*is a smooth immersion because $\gamma'(t)$ never vanishes. It is also injective, because $\gamma(t_1) = \gamma(t_2)$ implies that both $t_1 - t_2$ and $\alpha t_1 - \alpha t_2$ are integers, which is impossible unless $t_1 = t_2$.*

*Consider the set $\gamma(\mathbb{Z}) = \{\gamma(n) \mid n \in \mathbb{Z}\}$. It follows from Dirichlet's approximation theorem (see below) that for every $\varepsilon > 0$, there are integers $n, m$ such that $\|\alpha n - m\| < \varepsilon$. Using the fact that $\|e^{it_1} - e^{it_2}\| \le \|t_1 - t_2\|$ for $t_1, t_2 \in \mathbb{R}$ (because the line segment from $e^{it_1}$ to $e^{it_2}$ is shorter than the circular arc of length $\|t_1 - t_2\|$), we have*

$$\|e^{2\pi i \alpha n} - 1\| = \|e^{2\pi i \alpha n} - e^{2\pi i m}\| \le \|2\pi(\alpha n - m)\| < 2\pi\varepsilon.$$

*Therefore,*

$$\|\gamma(n) - \gamma(0)\| = \left\| \left( e^{2\pi i n}, e^{2\pi i \alpha n} \right) - (1,1) \right\| = \left\| \left( 1, e^{2\pi i \alpha n} \right) - (1,1) \right\| < 2\pi\varepsilon.$$

*Thus, $\gamma(0)$ is a limit point of $\gamma(\mathbb{Z})$. But this means that $\gamma$ is not a homeomorphism onto its image, because $\mathbb{Z}$ has no limit point in $\mathbb{R}$. In fact, it is not hard to show that the image set $\gamma(\mathbb{R})$ is actually dense in $\mathbb{T}^2$:*

*Let $(z^1, z^2) \in \mathbb{T}^2$ and choose $s^1, s^2 \in \mathbb{R}$ so that $e^{2\pi i s^1} = z^1$ and $e^{2\pi i s^2} = z^2$. Then*

$$\gamma(s^1 + n) = \left( e^{2\pi i s^1}, e^{2\pi i \alpha n} e^{2\pi i \alpha s^1} \right),$$

*so it remains to show that $S = \{e^{2\pi i \alpha n} \mid n \in \mathbb{Z}\}$ is dense in $\mathbb{S}^1$ whenever $\alpha$ is irrational. Let $\varepsilon > 0$; then by Dirichlet's approximation theorem, there exist integers $n, m$ such that $\|n\alpha - m\| < \varepsilon$. Let $\theta = n\alpha - m$ so that*

$$e^{2\pi i \theta} = e^{2\pi i \alpha n} \in S \qquad \text{and } \theta \ne 0 \text{ since } \alpha \text{ is irrational.}$$

*Any integer power of $e^{2\pi i \theta}$ is also a member of $S$, and since $\varepsilon$ was arbitrary, we can approximate any point on $\mathbb{S}^1$ by taking powers of $e^{2\pi i \theta}$, as desired.*                    ♟

The preceding example depends heavily on the following elementary result from number theory:

**Lemma 34 (Dirichlet's Approximation Theorem).** *Given $\alpha \in \mathbb{R}$ and any positive integer $N$, there exist integers $n, m$ with $1 \le n \le N$ such that $|n\alpha - m| < 1/N$.*

The following proposition gives a few simple sufficient criteria for an injective immersion to be an embedding:

**Proposition 63.** *Suppose M and N are smooth manifolds (with or without boundary), and $F\colon M \to N$ is an injective smooth immersion. If any of the following holds, then F is a smooth embedding.*

*a)* *F is an open or closed map.*

*b)* *F is a proper map.[2]*

*c)* *M is compact.*

*d)* *M has empty boundary and $\dim M = \dim N$.*

**Example 59.** *Let $\iota\colon \mathsf{S}^n \hookrightarrow \mathbb{R}^{n+1}$ be the inclusion map. We have previously shown that $\iota$ is smooth by computing its coordinate representation with respect to graph coordinates. It is easy to verify in the same coordinates that its differential is injective at each point, so it is an injective smooth immersion. Moreover, because $\mathsf{S}^n$ is compact, $\iota$ is a smooth embedding by part c) of the above proposition.* ❦

**Example 60.** *We now give an example of a smooth embedding that is neither an open nor a closed map. Let $X = [0,1)$ and $Y = [-1,1]$, and let $f\colon X \to Y$ be the identity map on X. Then f is a smooth embedding, X is both open and closed (in X), but $f(X)$ is neither open nor closed (in Y).* ❦

**Theorem 108** (**Local Embedding Theorem**). *Suppose M and N are smooth manifolds (with or without boundary), and $F\colon M \to N$ is a smooth map. Then F is a smooth immersion if and only if every point in M has a neighborhood $U \subseteq M$ such that $F|_U\colon U \to N$ is a smooth embedding.*

Theorem 108 points the way to a notion of immersions that makes sense for arbitrary topological spaces:

**Definition 149.** *If X and Y are topological spaces, a continuous map $F\colon X \to Y$ is called a **topological immersion** if every point of X has a neighborhood U such that $F|_U$ is a topological embedding.*

Thus, every smooth immersion is a topological immersion; but, just as with embeddings, a topological immersion that happens to be smooth need not be a smooth immersion (see Example 56 above).

✦────────────── ☞

[2] Recall that if X and Y are topological spaces, a map $F\colon X \to Y$ (continuous or not) is said to be ***proper*** if for every compact set $K \subseteq Y$, the preimage $F^{-1}(K)$ is compact as well.

### 3.5.3 Submersions

**Definition 150.** *If $\pi\colon M \to N$ is any continuous map, a **section of** $\pi$ is a continuous right inverse for $\pi$, i.e., a continuous map $\pi_S\colon N \to M$ such that $\pi \circ \pi_S = \mathrm{Id}_N$:*

$$
\begin{array}{c}
M \\
\pi \Big\downarrow \quad \Big) \pi_S \\
N
\end{array}
$$

*A **local section of** $\pi$ is a continuous map $\pi_S\colon U \to M$ defined on some open subset $U \subseteq N$ and satisfying the analogous relation $\pi \circ \pi_S = \mathrm{Id}_U$.*

**Theorem 109** (**Local Section Theorem**). *Suppose M and N are smooth manifolds and $\pi\colon M \to N$ is a smooth map. Then $\pi$ is a smooth submersion if and only if every point of M is in the image of a smooth local section of $\pi$.*

This theorem motivates the following definition:

**Definition 151.** *If $\pi\colon X \to Y$ is a continuous map, we say that $\pi$ is a **topological submersion** if every point of X is in the image of a (continuous) local section of $\pi$. (The preceding theorem shows that every smooth submersion is a topological submersion).*

**Example 61.** *For an example of a smooth map that is a topological submersion but not a smooth submersion, consider $f(x) = x^3$ at $x = 0$.*  ✈

**Proposition 64** (**Properties of Smooth Submersions**). *Let M and N be smooth manifolds, and suppose $\pi\colon M \to N$ is a smooth submersion. Then $\pi$ is an open map, and if it is surjective it is a quotient map.*

The next three theorems provide important tools that we will use frequently when studying submersions. They demonstrate that surjective smooth submersions play a role in smooth manifold theory analogous to the role of quotient maps in topology:

**Theorem 110 (Characteristic Property of Surjective Smooth Submersions).** *Suppose M and N are smooth manifolds, and $\pi\colon M \to N$ is a surjective smooth submersion. For any smooth manifold S (with or without boundary), a map $F\colon N \to S$ is smooth if and only if $F \circ \pi$ is smooth:*

$$
\begin{array}{ccc}
M & & \\
\Big\downarrow{\scriptstyle\pi} & \searrow{\scriptstyle F \circ \pi} & \\
N & \xrightarrow{\ \ F\ \ } & S
\end{array}
$$

**Side Note from Topology:** Recall that if $\pi\colon X \to Y$ is a map, a subset $U \subseteq X$ is said to be **saturated with respect to $\pi$** if $U$ is the entire preimage of its image under $\pi$, i.e., if $U = \pi^{-1}\left(\pi(U)\right)$. Given $y \in Y$, the **fiber of $\pi$ over $y$** is the set $\pi^{-1}(y)$. Thus, a subset of $X$ is saturated if and only if it is a union of fibers.

**Theorem 111 (Passing Smoothly to the Quotient).** *Suppose M and N are smooth manifolds and $\pi\colon M \to N$ is a surjective smooth submersion. If S is a smooth manifold (with or without boundary) and $F\colon M \to S$ is a smooth map that is constant on the fibers of $\pi$, then there exists a unique smooth map $\widetilde{F}\colon N \to S$ such that $\widetilde{F} \circ \pi = F$:*

$$
\begin{array}{ccc}
M & & \\
\Big\downarrow{\scriptstyle\pi} & \searrow{\scriptstyle F} & \\
N & \dashrightarrow{\scriptstyle \widetilde{F}} & S
\end{array}
$$

**Theorem 112 (Uniqueness of Smooth Quotients).** *Suppose that $M$, $N_1$, and $N_2$ are smooth manifolds, and $\pi_1\colon M \to N_1$ and $\pi_2\colon M \to N_2$ are surjective smooth submersions that are constant on each other's fibers. Then there exists a unique diffeomorphism $F\colon N_1 \to N_2$ such that $F \circ \pi_1 = \pi_2$:*

$$
\begin{array}{ccc}
 & M & \\
{\scriptstyle\pi_1}\swarrow & & \searrow{\scriptstyle\pi_2} \\
N_1 & \dashrightarrow{\scriptstyle F} & N_2
\end{array}
$$

### 3.5.4   Smooth Covering Maps

Recall from topology the notion of covering maps (for more on this subject, check Subsection §**??**):

**Definition 152.** *Suppose E and X are topological spaces. A map* $\pi\colon E \to X$ *is called a **covering map** if E and X are connected and locally path-connected, $\pi$ is surjective and continuous, and each point $p \in X$ has a neighborhood U that is **evenly covered by** $\pi$, meaning that each component of $\pi^{-1}(U)$ is mapped homeomorphically onto U by $\pi$. In this case, X is called the **base of the covering**, and E is called a **covering space of** X. If U is an evenly covered subset of X, the components of $\pi^{-1}(U)$ are called the **sheets of the covering over** U.*

In the context of smooth manifolds, it is useful to introduce a slightly more restrictive type of covering map:

**Definition 153.** *If E and M are connected[3] smooth manifolds (with or without boundary), a map* $\pi\colon E \to M$ *is called a **smooth covering map** if $\pi$ is smooth and surjective, and each point in M has a neighborhood U such that each component of $\pi^{-1}(U)$ is mapped diffeomorphically onto U by $\pi$. In this context we also say that U is **evenly covered by** $\pi$ and that the space M is called the **base of the covering**, while E is called a **covering manifold of** M. If E is simply connected, it is called the **universal covering manifold of** M.*

   **Remark:** To distinguish this new definition from the previous one, we often call an ordinary (not necessarily smooth) covering map a ***topological covering map***. A smooth covering map is, in particular, a topological covering map. But as with other types of maps we have studied in this chapter, a smooth covering map is more than just a topological covering map that happens to be smooth: the definition requires in addition that the restriction of $\pi$ to each component of the preimage of an evenly covered set be a diffeomorphism, not just a smooth homeomorphism.

**Proposition 65 (Properties of Smooth Coverings).** *We have the following properties of smooth coverings:*

---

[3] Note that here we make no mention of path-connectedness, since, as you may well recall, connectedness $\implies$ path-connectedness in manifolds.

*a)* *Every smooth covering map is a local diffeomorphism, a smooth submersion, an open map, and a quotient map.*

*b)* *An injective smooth covering map is a diffeomorphism.*

*c)* *A topological covering map is a smooth covering map if and only if it is a local diffeomorphism.*

Because smooth covering maps are surjective smooth submersions, all of the previous results about smooth submersions can be applied to them. For example, Theorem 111 is a particularly useful tool for defining a smooth map out of the base of a covering space.

For smooth covering maps, the local section theorem can be strengthened:

**Proposition 66** (**Local Section Theorem for Smooth Covering Maps**). *Suppose E and M are smooth manifolds (with or without boundary), and $\pi: E \to M$ is a smooth covering map. Given any evenly covered open subset $U \subseteq M$, any $q \in U$, and any p in the fiber of $\pi$ over q, there exists a unique smooth local section $\pi_S: U \to E$ such that $\pi_S(q) = p$.*

**Proposition 67** (**Covering Spaces of Smooth Manifolds**). *Suppose M is a connected smooth n-manifold, and $\pi: E \to M$ is a topological covering map. Then E is a topological n-manifold, and has a unique smooth structure such that $\pi$ is a smooth covering map.*

**Proposition 68** (**Covering Spaces of Smooth Manifolds with Boundary**). *Suppose M is a connected smooth n-manifold with boundary, and $\pi: E \to M$ is a topological covering map. Then E is a topological n-manifold with boundary such that $\partial E = \pi^{-1}(\partial M)$, and it has a unique smooth structure such that $\pi$ is a smooth covering map.*

**Corollary 32** (**Existence of a Universal Covering Manifold**). *If M is a connected smooth manifold, there exists a simply connected smooth manifold $\widetilde{M}$, called the **universal covering manifold of** M, and a smooth covering map $\pi: \widetilde{M} \to M$. The universal covering manifold is unique in the following sense: if $\widetilde{M}'$ is any other simply connected smooth manifold that admits a smooth covering map $\pi': \widetilde{M}' \to M$, then there exists a diffeomorphism $\Phi: \widetilde{M} \to \widetilde{M}'$ such that $\pi' \circ \Phi = \pi$.*

There are not many simple criteria for determining whether a given map is a smooth covering map, even if it is known to be a surjective local diffeomorphism. The following

proposition gives one useful sufficient criterion. (It is not a necessary condition, however; see Problem $4 - 11$ on [Lee, 2013, p. 96].)

**Proposition 69** (**Covering Spaces of Smooth Manifolds with Boundary**)**.** *Suppose E and M are nonempty connected smooth manifolds (with or without boundary). If $\pi\colon E \to M$ is a proper local diffeomorphism, then $\pi$ is a smooth covering map.*

## 3.6 Submanifolds

In this section we will explore smooth submanifolds, which are smooth manifolds on their own right, that happen to be subsets of other smooth manifolds. As you will soon discover, the situation is quite a bit more subtle and complex than the analogous theory of topological spaces/subspaces. I should mention that, even though we will not discuss it in these notes, there is a very celebrated result known as *Whitney's Embedding Theorem* that you should definitely check out after you're done with this section. You can read it on [Lee, 2013, Chapter 6], or just look it up somewhere online. The theorem basically states that abstract manifolds are not so abstract after all! It turns about that every manifold can be embedded into some ambient Euclidean space! This is a beautiful theorem that you should definitely investigate, but first you need to make sure that you thoroughly understand the content on this section so hang on cowboy, not so fast!

### 3.6.1 Embedded Submanifolds

**Definition 154.** *Suppose M is a smooth manifold (with or without boundary). An **embedded submanifold** (also called **regular submanifold**) of M is a subset $S \subseteq M$ that is a manifold (without boundary) in the subspace topology, endowed with a smooth structure with respect to which the inclusion map $S \hookrightarrow M$ is a smooth embedding.*

**Definition 155.** *If S is an embedded submanifold of M, the difference* dim *M* − dim *S is called the* **codimension of** *S* **in** *M, and the containing manifold M is called the* **ambient manifold for** *S. An* **embedded hypersurface** *is an embedded submanifold of codimension* 1. *(The empty set is an embedded submanifold of any dimension).*

The easiest embedded submanifolds to understand are those of codimension 0. Recall that for any smooth manifold *M* we defined an open submanifold of *M* to be any open subset with the subspace topology and with the smooth charts obtained by restricting those of *M*.

**Proposition 70 (Open Submanifolds).** *Suppose M is a smooth manifold. The embedded submanifolds of codimension* 0 *in M are exactly the open submanifolds.*

The next few propositions demonstrate several other ways to produce embedded submanifolds:

**Proposition 71 (Images of Embeddings as Submanifolds).** *Suppose M is a smooth manifold (with or without boundary), N is a smooth manifold, and F : N → M is a smooth embedding. Let S = F(N). With the subspace topology, S is a topological manifold, and it has a unique smooth structure making it into an embedded submanifold of M with the property that F is a diffeomorphism onto its image.*

*Proof.* If we give *S* the subspace topology that it inherits from *M*, then the assumption that *F* is an embedding means that *F* can be considered as a homeomorphism from *N* onto *S*, and thus *S* is a topological manifold. We give *S* a smooth structure by taking the smooth charts to be those of the form $(F(U), \varphi \circ F^{-1})$, where $(U, \varphi)$ is any smooth chart for *N*; smooth compatibility of these charts follows immediately from the smooth compatibility of the corresponding charts for *N*. With this smooth structure on *S*, the map *F* is a diffeomorphism onto its image (essentially by definition), and this is obviously the only smooth structure with this property. The inclusion map *S* ↪ *M* is equal to the composition of a diffeomorphism followed by a smooth embedding:

$$S \xrightarrow{F^{-1}} N \xrightarrow{F} M,$$

and therefore it is a smooth embedding.                                                  □

Since every embedded submanifold is the image of a smooth embedding (namely its own inclusion map), the previous proposition shows that embedded submanifolds are exactly the images of smooth embeddings.

**Proposition 72 (Slices of Product Manifolds).** *Suppose M and N are smooth manifolds. For each $p \in N$, the subset $M \times \{p\}$ (called a **slice** of the product manifold) is an embedded submanifold of $M \times N$ diffeomorphic to M.*

*Proof.* The set $M \times \{p\}$ is the image of the smooth embedding $x \mapsto (x, p)$.                □

**Proposition 73 (Graphs as Submanifolds).** *Suppose M is a smooth m-manifold (without boundary), N is a smooth n-manifold (with or without boundary), $U \subseteq M$ is open, and $f : U \to N$ is a smooth map. Let $\Gamma(f) \subseteq M \times N$ denote the graph of $f$:*

$$\Gamma(f) = \{(x, y) \in M \times N \mid x \in U, \ y = f(x)\}.$$

*Then $\Gamma(f)$ is an embedded m-dimensional submanifold of $M \times N$ (see* Figure 3.12*).*



Figure 3.12: A graph is an embedded submanifold.

*Proof.* Define a map $\gamma_f : U \to M \times N$ by

$$\gamma_f(x) = (x, f(x)). \tag{3.4}$$

It is a smooth map whose image is $\Gamma(f)$. Because the projection $\pi_M : M \times N \to M$ satisfies $\pi_M \circ \gamma_f(x) = x$ for $x \in U$, the composition $d(\pi_M)_{(x, f(x))} \circ d(\gamma_f)_x$ is the identity on $T_x M$ for each $x \in U$. Thus, $d(\gamma_f)_x$ is injective, so $f$ is a smooth immersion. It a homeomorphism onto its image because $\pi_M|_{\Gamma(f)}$ is a continuous inverse for it. Thus, $\Gamma(f)$ is an embedded submanifold diffeomorphic to $U$.                □

Recall the following proposition from topology:

**Proposition 74 (Sufficient Conditions for Properness).** *Suppose X and Y are topological spaces, and $F\colon X \to Y$ is a continuous map.*

  *a)* *If X is compact and Y is Hausdorff, then F is proper.*

  *b)* *If F is a closed map with compact fibers, then F is proper.*

  *c)* *If F is a topological embedding with closed image, then F is proper.*

  *d)* *If Y is Hausdorff and F has a continuous left inverse (i.e., a continuous map $G\colon Y \to X$ such that $G \circ F = \mathrm{Id}_X$), then F is proper.*

  *e)* *If F is proper and $A \subseteq X$ is a subset that is saturated with respect to F, then $F|_A\colon A \to F(A)$ is proper.*

**Definition 156.** *An embedded submanifold $S \subseteq M$ is said to be **properly embedded** if the inclusion $S \hookrightarrow M$ is a proper map.*

**Proposition 75.** *Suppose M is a smooth manifold (with or without boundary) and $S \subseteq M$ is an embedded submanifold. Then S is properly embedded if and only if it is a closed subset of M.*

**Corollary 33.** *Every compact embedded submanifold is properly embedded.*

*Proof.* Compact subsets of Hausdorff spaces are closed.                                    □

   Graphs of globally defined functions are common examples of properly embedded submanifolds:

**Proposition 76 (Global Graphs Are Properly Embedded).** *Suppose M is a smooth manifold, N is a smooth manifold (with or without boundary), and $f\colon M \to N$ is a smooth map. With the smooth manifold structure of* Proposition 73, $\Gamma(f)$ *is properly embedded in $M \times N$.*

*Proof.* Proof. In this case, the projection $\pi_M\colon M \times N \to M$ is a smooth left inverse for the embedding $\gamma_f\colon M \to M \times N$ defined by equation (3.4) above. Thus $\gamma_f$ is proper by Proposition 74.                                    □

As Theorem 113 below will show, embedded submanifolds are modeled locally on the standard embedding of $\mathbb{R}^k$ into $\mathbb{R}^n$, identifying $\mathbb{R}^k$ with the subspace

$$\{(x^1, \ldots, x^k, x^{k+1}, \ldots, x^n) \mid x^{k+1} = \cdots = x^n = 0\} \subseteq \mathbb{R}^n.$$

Somewhat more generally:

**Definition 157.** *If $U$ is an open subset of $\mathbb{R}^n$ and $k \in \{0, \ldots, n\}$, a k-**dimensional slice of** $U$ (or simply a k-**slice**) is any subset of the form*

$$S = \{(x^1, \ldots, x^k, x^{k+1}, \ldots, x^n) \in U \mid x^{k+1} = c^{k+1}, \ldots, x^n = c^n\}$$

*for some constants $c^{k+1}, \ldots, c^n$. (When $k = n$, this just means $S = U$.)*

**Definition 158.** *Let $M$ be a smooth n-manifold, and let $(U, \varphi)$ be a smooth chart on $M$. If $S$ is a subset of $U$ such that $\varphi(S)$ is a k-slice of $\varphi(U)$, then we say that $S$ is a k-**slice of** $U$.*

**Remark:** Although in general we allow our slices to be defined by arbitrary constants $c^{k+1}, \ldots, c^n$, it is sometimes useful to have slice coordinates for which the constants are all zero, which can easily be achieved by subtracting a constant from each coordinate function.

**Definition 159.** *Given a subset $S \subseteq M$ and a nonnegative integer $k$, we say that $S$ satisfies the **local k-slice condition** if each point of $S$ is contained in the domain of a smooth chart $(U, \varphi)$ for $M$ such that $S \cap U$ is a single k-slice in $U$. Any such chart is called a **slice chart for** $S$ **in** $M$, and the corresponding coordinates $(x^1, \ldots, x^n)$ are called **slice coordinates**.*

**Theorem 113** (**Local Slice Criterion for Embedded Submanifolds**). *Let $M$ be a smooth n-manifold. If $S \subseteq M$ is an embedded k-dimensional submanifold, then $S$ satisfies the local k-slice condition. Conversely, if $S \subseteq M$ is a subset that satisfies the local k-slice condition, then with the subspace topology, $S$ is a topological manifold of dimension k, and it has a smooth structure making it into a k-dimensional embedded submanifold of M.*

*Proof.* $(\Rightarrow)$ First suppose that $S \subseteq M$ is an embedded $k$-dimensional submanifold. Since the inclusion map $S \hookrightarrow M$ is an immersion, the rank theorem shows that for any $p \in S$

there are smooth charts $(U, \varphi)$ for $S$ (in its given smooth manifold structure) and $(V, \psi)$ for $M$, both centered at $p$, in which the inclusion map $\iota|_U \colon U \to V$ has the coordinate representation

$$(x^1, \ldots, x^k) \mapsto (x^1, \ldots, x^k, 0, \ldots, 0).$$

Choose $\varepsilon > 0$ small enough that both $U$ and $V$ contain coordinate balls of radius $\varepsilon$ centered at $p$, and denote these coordinate balls by $U_0 \subseteq U$ and $V_0 \subseteq V$. It follows that $U_0 = \iota(U_0)$ is exactly a single slice in $V_0$. Because $S$ has the subspace topology, the fact that $U_0$ is open in $S$ means that there is an open subset $W \subseteq M$ such that $U_0 = W \cap S$. Setting $V_1 = V_0 \cap W$, we obtain a smooth chart $\left(V_1, \psi|_{V_1}\right)$ for $M$ containing $p$ such that $V_1 \cap S = U_0$, which is a single slice of $V_1$.

($\Leftarrow$) Conversely, suppose $S$ satisfies the local $k$-slice condition. With the subspace topology, $S$ is Hausdorff and second countable, because both properties are inherited by subspaces. To see that $S$ is locally Euclidean, we construct an atlas. The basic idea of the construction is that if $(x^1, \ldots, x^n)$ are slice coordinates for $S$ in $M$, we can use $(x^1, \ldots, x^k)$ as local coordinates for $S$.

For this proof, let $\pi \colon \mathbb{R}^n \to \mathbb{R}^k$ denote the projection onto the first $k$ coordinates. Let $(U, \varphi)$ be any slice chart for $S$ in $M$ (see Figure 3.13), and define

$$V = U \cap S, \qquad \widehat{V} = \pi \circ \varphi(V), \qquad \psi = \pi \circ \varphi|_V \colon V \to \widehat{V}.$$



Figure 3.13: A chart for a subset satisfying the k-slice condition.

By definition of slice charts, $\varphi(V)$ is the intersection of $\varphi(U)$ with a certain $k$-slice $A \subseteq \mathbb{R}^n$ defined by setting $x^{k+1} = c^{k+1}, \ldots, x^n = c^n$, and therefore $\varphi(V)$ is open in $A$. Since $\pi|_A$ is a diffeomorphism from $A$ to $\mathbb{R}^k$, it follows that $\widehat{V}$ is open in $\mathbb{R}^k$. Moreover, $\psi$ is a

homeomorphism because it has a continuous inverse given by $\varphi^{-1} \circ j|_{\widehat{V}}$, where $j \colon \mathbb{R}^k \to \mathbb{R}^n$ is the map

$$j(x^1, \ldots, x^k) = (x^1, \ldots, x^k, c^{k+1}, \ldots, c^n)$$

Thus $S$ is a topological $k$-manifold, and the inclusion map $\iota \colon S \hookrightarrow M$ is a topological embedding.

To put a smooth structure on $S$, we need to verify that the charts constructed above are smoothly compatible. Suppose $(U, \varphi)$ and $(U', \varphi')$ are two slice charts for $S$ in $M$, and let $(V, \psi)$, $(V', \psi')$ be the corresponding charts for $S$. The transition map is given by $\psi' \circ \psi^{-1} = \pi \circ \varphi' \circ \varphi^{-1} \circ j$, which is a composition of four smooth maps (see Figure 3.14)



Figure 3.14: Smooth compatibility of slice charts.

Thus the atlas we have constructed is in fact a smooth atlas, and it defines a smooth structure on $S$. In terms of a slice chart $(U, \varphi)$ for $M$ and the corresponding chart $(V, \psi)$ for $S$, the inclusion map $S \hookrightarrow M$ has a coordinate representation of the form

$$(x^1, \ldots, x^k) \mapsto (x^1, \ldots, x^k, c^{k+1}, \ldots, c^n),$$

which is a smooth immersion. Since the inclusion is a smooth immersion and a topological embedding, $S$ is an embedded submanifold. $\qquad \square$

Notice that the local slice condition for $S \subseteq M$ is a condition on the subset $S$ only; it does not presuppose any particular topology or smooth structure on $S$. As we will see later on on Theorem 118, the smooth manifold structure constructed in the preceding theorem is the unique one in which $S$ can be considered as a submanifold, so a subset satisfying the local slice condition is an embedded submanifold in only one way.

**Example 62 (Spheres as Submanifolds).** *For any $n \geq 0$, $\mathbb{S}^n$ is an embedded submanifold of $\mathbb{R}^{n+1}$, because it is locally the graph of a smooth function: the intersection of $\mathbb{S}^n$ with the open subset $\{x \colon x^i > 0\}$ is the graph of the smooth function*

$$x^i = f(x^1, \ldots, x^{i-1}, x^{i+1}, \ldots, x^{n+1}),$$

*where $f \colon \mathbb{B}^n \to \mathbb{R}$ is given by $f(u) = \sqrt{1 - \|u\|^2}$. Similarly, the intersection of $\mathbb{S}^n$ with $\{x \mid x^i < 0\}$ is the graph of $-f$. Since every point in $\mathbb{S}^n$ is in one of these sets, $\mathbb{S}^n$ satisfies the local n-slice condition and is thus an embedded submanifold of $\mathbb{R}^{n+1}$. The smooth structure thus induced on $\mathbb{S}^n$ is the same as the one we have previously defined in class: in fact, the coordinates for $\mathbb{S}^n$ determined by these slice charts are exactly the graph coordinates previously defined.* ✦

If $M$ is a smooth manifold with nonempty boundary and $S \subseteq M$ is an embedded submanifold, then $S$ might intersect $\partial M$ in very complicated ways, so we will not attempt to prove any general results about the existence of slice charts for $S$ in $M$ in that case. However, in the special case in which the submanifold is the boundary of $M$ itself, the boundary charts for $M$ play the role of slice charts for $\partial M$ in $M$, and we do have the following result:

**Theorem 114.** *If $M$ is a smooth n-manifold with boundary, then with the subspace topology, $\partial M$ is a topological $(n-1)$-dimensional manifold (without boundary), and has a smooth structure such that it is a properly embedded submanifold of $M$.*

### 3.6.2    Level Sets

**Definition 160.** *If $\Phi \colon M \to N$ is any map and $c$ is any point of $N$, we call the set $\Phi^{-1}(c)$ a **level set of $\Phi$** (see Figure 3.15). (In the special case when $N = \mathbb{R}^k$ and $c = 0$, the level set $\Phi^{-1}(0)$ is usually called the **zero set of $\Phi$**).*

**Example 63.** *It is easy to find level sets of smooth functions that are not smooth submanifolds. For instance, consider the three smooth functions $\Theta, \Phi, \Psi \colon \mathbb{R}^2 \to \mathbb{R}$ defined by*

$$\Theta(x, y) = x^2 - y, \qquad \Phi(x, y) = x^2 - y^2, \qquad \Psi(x, y) = x^2 - y^3.$$

*Although the zero set of $\Theta$ (a parabola) is an embedded submanifold of $\mathbb{R}^2$ (because it is the graph of the smooth function $f(x) = x^2$), it can be shown that neither the zero set of $\Phi$ nor that of*

Figure 3.15: A level set.



Figure 3.16: Level sets may or may not be embedded submanifolds.

*Ψ is an embedded submanifold (I'll leave this to you as an exercise; try solving* Problem 5-11 *from [Lee, 2013]). In fact, without further assumptions on the smooth function, the situation is about as bad as could be imagined: as it has been previously stated on Theorem 104 in Section §3.3, every closed subset of M can be expressed as the zero set of some smooth real-valued function.* 🌍

**Theorem 115 (Constant-Rank Level Set Theorem).** *Let M and N be smooth manifolds, and let Φ: M → N be a smooth map with constant rank r. Each level set of Φ is a properly embedded submanifold of codimension r in M.*

**Corollary 34 (Submersion Level Set Theorem).** *If M and N be smooth manifolds and Φ: M → N is a smooth submersion, then each level set of Φ is a properly embedded submanifold whose codimension is equal to the dimension of N.*

*Proof.* The proof follows trivially by the fact that every smooth submersion has constant rank equal to the dimension of its codomain. □

This last result should be compared to the corresponding result in linear algebra: if $L\colon \mathbb{R}^m \to \mathbb{R}^r$ is a surjective linear map, then the kernel of $L$ is a linear subspace of codimension $r$ by the rank-nullity law. The vector equation $Lx = 0$ is equivalent to $r$ linearly independent scalar equations, each of which can be thought of as cutting down one of the degrees of freedom in $\mathbb{R}^m$, leaving a subspace of codimension $r$. In the context of smooth manifolds, the analogue of a surjective linear map is a smooth submersion, each of whose (local) component functions cuts down the dimension by one.

Corollary 34 can be strengthened considerably, because we need only check the submersion condition on the level set we are interested in.

**Definition 161.** *If $\Phi\colon M \to N$ is a smooth map, a point $p \in M$ is said to be a **regular point of $\Phi$** if $\mathrm{d}\Phi_p\colon T_pM \to T_{\Phi(p)}N$ is surjective; it is a **critical point of $\Phi$** otherwise. (This means, in particular, that every point of $M$ is critical if $\dim M < \dim N$, and every point is regular if and only if $F$ is a submersion.)*

**Definition 162.** *A point $c \in N$ is said to be a **regular value of $\Phi$** if every point of the level set $\Phi^{-1}(c)$ is a regular point; $c$ is called a **critical value** otherwise. (In particular, if $\Phi^{-1}(c) = \varnothing$, then $c$ is a regular value.*

**Definition 163.** *A level set $\Phi^{-1}(c)$ is called a **regular level set** if $c$ is a regular value of $\Phi$; in other words, a regular level set is a level set consisting entirely of regular points of $\Phi$ (points $p$ such that $\mathrm{d}\Phi_p$ is surjective).*

**Corollary 35** (**Regular Level Set Theorem**)**.** *Every regular level set of a smooth map between smooth manifolds is a properly embedded submanifold whose codimension is equal to the dimension of the codomain.*

It is worth noting that the previous corollary also applies to empty level sets, which are both regular level sets and properly embedded submanifolds.

**Example 64** (**Spheres**)**.** *Now we can give a much easier proof that $\mathbb{S}^n$ is an embedded submanifold of $\mathbb{R}^{n+1}$. The sphere is a regular level set of the smooth function $f\colon \mathbb{R}^{n+1} \to \mathbb{R}$ given by $f(x) = \|x\|^2$, since $\mathrm{d}f_x(v) = 2\sum_i x^i v^i$ which is surjective (except at the origin, but obviously the origin is not a point on the sphere).*

Not all embedded submanifolds can be expressed as level sets of smooth submersions. However, the next proposition shows that every embedded submanifold is at least locally of this form:

**Proposition 77.** *Let S be a subset of a smooth m-manifold M. Then S is an embedded k-submanifold of M if and only if every point of S has a neighborhood U in M such that $U \cap S$ is a level set of a smooth submersion $\Phi: U \to \mathbb{R}^{m-k}$.*

*Proof.* ($\Rightarrow$) First suppose $S$ is an embedded $k$-submanifold. If $(x^1, \ldots, x^m)$ are slice coordinates for $S$ on an open subset $U \subseteq M$, then the map $\Phi: U \to \mathbb{R}^{m-k}$ given in coordinates by $\Phi(x) = (x^{k+1}, \ldots, x^m)$ is easily seen to be a smooth submersion, one of whose level sets is $S \cap U$ (see Figure 3.17).



Figure 3.17: An embedded submanifold is locally a level set.

($\Leftarrow$) Conversely, suppose that around every point $p \in S$ there is a neighborhood $U$ and a smooth submersion $\Phi: U \to \mathbb{R}^{m-k}$ such that $S \cap U$ is a level set of $\Phi$. Then by the submersion level set theorem (Corollary 34), we have that $S \cap U$ is an embedded submanifold of $U$, so it satisfies the local slice condition; it follows that $S$ is itself an embedded submanifold of $M$. $\qquad\square$

**Definition 164.** *If $S \subseteq M$ is an embedded submanifold, a smooth map $\Phi: M \to N$ such that $S$ is a regular level set of $\Phi$ is called a **defining map for $S$**. In the special case when $N = \mathbb{R}^{m-k}$ (so that $\Phi$ is a real-valued or vector-valued function), it is usually called a **defining function**.*

For instance, from Example 64, we have that $f(x) = \|x\|^2$ is a defining function for the sphere. More generally:

**Definition 165.** *If U is an open subset of M and $\Phi\colon U \to N$ is a smooth map such that $S \cap U$ is a regular level set of $\Phi$, then $\Phi$ is called a **local defining map** (or **local defining function**) **for S**.*

Proposition 77 says that every embedded submanifold admits a local defining function in a neighborhood of each of its points. In specific examples, finding a (local or global) defining function for a submanifold is usually just a matter of using geometric information about how the submanifold is defined together with some computational ingenuity. Here is an example:

**Example 65 (Surfaces of Revolution).** *Let H be the half-plane $\{(r,z) \mid r > 0\}$, and suppose $C \subseteq H$ is an embedded $1$-dimensional submanifold. The surface of revolution determined by C is the subset $S_C \subseteq \mathbb{R}^3$ given by*

$$S_C = \left\{ (x,y,z) : \left( \sqrt{x^2 + y^2}, z \right) \in C \right\}.$$

*The set C is called its **generating curve** (see Figure 3.18).*



Figure 3.18: A surface of revolution.

*If $\varphi\colon U \to \mathbb{R}$ is any local defining function for C in H, we get a local defining function $\Phi$ for $S_C$ by*

$$\Phi(x,y,z) = \varphi\left( \sqrt{x^2 + y^2}, z \right),$$

*defined on the open subset*

$$\widetilde{U} = \left\{ (x,y,z) : \left( \sqrt{x^2 + y^2}, z \right) \in U \right\} \subseteq \mathbb{R}^3.$$

*A computation shows that the Jacobian matrix of $\Phi$ is*

$$D\Phi(x,y,z) = \left( \frac{x}{r} \frac{\partial \varphi}{\partial r}(r,z) \quad \frac{y}{r} \frac{\partial \varphi}{\partial r}(r,z) \quad \frac{\partial \varphi}{\partial z}(r,z) \right),$$

*where we have written $r = \sqrt{x^2 + y^2}$. At any point $(x,y,z) \in S_C$, at least one of the components of $D\Phi(x,y,z)$ is nonzero, so $S_C$ is a regular level set of $\Phi$ and is thus an embedded 2-dimensional submanifold of $\mathbb{R}^3$.*

*For a specific example, the doughnut-shaped torus of revolution is the surface of revolution obtained from the circle $(r - 2)^2 + z^2 = 1$. It is a regular level set of the function $\Phi(x,y,z) = \left( \sqrt{x^2 + y^2} - 2 \right)^2 + z^2$, which is smooth on $\mathbb{R}^3$ minus the z-axis.*  ⬤

### 3.6.3   Immersed Submanifolds

**Definition 166.** *Suppose M is a smooth manifold (with or without boundary). An **immersed submanifold of M** is a subset $S \subseteq M$ endowed with a topology (not necessarily the subspace topology) with respect to which it is a topological manifold (without boundary), and a smooth structure with respect to which the inclusion map $S \hookrightarrow M$ is a smooth immersion. (As in the case of embedded submanifolds, we define the codimension of S in M to be $\dim M - \dim S$).*

Note that every embedded submanifold is also an immersed submanifold. Because immersed submanifolds are the more general of the two types of submanifolds, we adopt the convention that the term *smooth submanifold* without further qualification means an immersed one, which includes an embedded submanifold as a special case. Similarly, the term *smooth hypersurface* without qualification means an immersed submanifold of codimension 1.

Immersed submanifolds often arise in the following way:

**Proposition 78 (Images of Immersions as Submanifolds).** *Suppose M is a smooth manifold (with or without boundary), N is a smooth manifold, and $F\colon N \to M$ is an injective smooth immersion. Let $S = F(N)$. Then S has a unique topology and smooth structure such that it is a smooth submanifold of M and such that $F\colon N \to S$ is a diffeomorphism onto its image.*

The following observation is sometimes useful when thinking about the topology of an immersed submanifold:

**Proposition 79.** *Suppose M is a smooth manifold and $S \subseteq M$ is an immersed submanifold. Then every subset of S that is open in the subspace topology is also open in its given submanifold topology, and the converse is true if and only if S is embedded.*

**Proposition 80 (Sufficient Conditions for Immersed Submanifolds to be Embedded).** *Suppose M is a smooth manifold (with or without boundary), and $S \subseteq M$ is an immersed submanifold. If any of the following holds, then S is embedded.*

   *a)* *S has codimension 0 in M.*

   *b)* *The inclusion map $S \hookrightarrow M$ is proper.*

   *c)* *S is compact.*

Although many immersed submanifolds are not embedded, the next proposition shows that the local structure of an immersed submanifold is the same as that of an embedded one:

**Proposition 81 (Immersed Submanifolds Are Locally Embedded).** *If M is a smooth manifold (with or without boundary), and $S \subseteq M$ is an immersed submanifold, then for each $p \in S$ there exists a neighborhood U of p in S that is an embedded submanifold of M.*

*Proof.* By Theorem 108 from Section §3.5, we have that each $p \in S$ has a neighborhood $U$ in $S$ such that the inclusion $\iota|_U \colon U \hookrightarrow M$ is an embedding. $\qquad\square$

**Remark:** It is important to be clear about what this proposition does and does not say: given an immersed submanifold $S \subseteq M$ and a point $p \in S$, it is possible to find a neighborhood $U$ of $p$ (in $S$) such that $U$ is embedded; but it may not be possible to find a neighborhood $V$ of $p$ (in $M$) such that $V \cap S$ is embedded (see Figure 3.19).

Figure 3.19: An immersed submanifold is locally embedded.

**Definition 167.** *Suppose $S \subseteq M$ is an immersed k-dimensional submanifold. A **local parametriza-tion** of S is a continuous map $\Psi: U \to M$ whose domain is an open subset $U \subseteq \mathbb{R}^k$, whose image is an open subset of S, and which, considered as a map into S, is a homeomorphism onto its image. It is called a **smooth local parametrization** if it is a diffeomorphism onto its image (with respect to the smooth manifold structure of S). If the image of $\Psi$ is all of S, it is called a **global parametrization**.*

**Proposition 82.** *Suppose M is a smooth manifold (with or without boundary), $S \subseteq M$ is an immersed k-submanifold, $\iota: S \hookrightarrow M$ is the inclusion map, and U is an open subset of $\mathbb{R}^k$. A map $\Psi: U \to M$ is a smooth local parametrization of S if and only if there is a smooth coordinate chart $(V, \varphi)$ for S such that $\Psi = \iota \circ \varphi^{-1}$. Therefore, every point of S is in the image of some local parametrization.*

### 3.6.4   Restricting Maps to Submanifolds

Given a smooth map $F: M \to N$, it is important to know whether $F$ is still smooth when its domain or codomain is restricted to a submanifold. In the case of restricting the domain, the answer is easy:

**Theorem 116 (Restricting the Domain of a Smooth Map).** *If M and N are smooth manifolds (with or without boundary), $F: M \to N$ is a smooth map, and $S \subseteq M$ is an immersed or embedded submanifold (see Figure 3.20), then $F\big|_S: S \to N$ is smooth.*

*Proof.* The inclusion map $\iota: S \hookrightarrow M$ is smooth by definition of an immersed submanifold. Since $F\big|_S = F \circ \iota$, the result follows.                                                    $\square$

Figure 3.20: Restricting the domain.

When the codomain is restricted, however, the resulting map may not be smooth, as the following example shows:

**Example 66.** *Let $S \subseteq \mathbb{R}^2$ be the figure-eight submanifold (lemniscate), with the topology and smooth structure induced by the immersion $\beta\colon (-\pi, \pi) \to \mathbb{R}^2$ defined by $\beta(t) = (\sin 2t, \sin t)$. Define a smooth map $G\colon \mathbb{R} \to \mathbb{R}^2$ by $G(t) = (\sin 2t, \sin t)$. (We are using the same formula that we used to define $\beta$, but now the domain is extended to the whole real line instead of being just a subinterval.) It is easy to check that the image of $G$ lies in $S$. However, as a map from $\mathbb{R}$ to $S$, $G$ is not even continuous, because $\beta^{-1} \circ G$ is not continuous at $t = \pi$.* ✈

The next theorem gives sufficient conditions for a map to be smooth when its codomain is restricted to an immersed submanifold. It shows that the failure of continuity is the only thing that can go wrong:

**Theorem 117 (Restricting the Codomain of a Smooth Map).** *Suppose $M$ is a smooth manifold (without boundary), $S \subseteq M$ is an immersed submanifold, and $F\colon N \to M$ is a smooth map whose image is contained in $S$ (see Figure 3.21). If $F$ is continuous as a map from $N$ to $S$, then $F\colon N \to S$ is smooth.*

This theorem is stated only for the case in which the ambient manifold $M$ is a manifold without boundary, because it is only in that case that we have constructed slice charts for embedded submanifolds of $M$. But the conclusion of the theorem is still true when $M$ has nonempty boundary (see Problem $9 - 13$ from [Lee, 2013]).

In the special case in which the submanifold $S$ is embedded, the continuity hypothesis is always satisfied:

Figure 3.21: Restricting the codomain.

**Corollary 36** (**Embedded Case**)**.** *Let $M$ be a smooth manifold and $S \subseteq M$ be an embedded submanifold. Then every smooth map $F\colon N \to M$ whose image is contained in $S$ is also smooth as a map from $N$ to $S$.*

*Proof.* Since $S \subseteq M$ has the subspace topology, a continuous map $F\colon N \to M$ whose image is contained in $S$ is automatically continuous into $S$, by the characteristic property of the subspace topology. □

Although the conclusion of the preceding corollary fails for some immersed submanifolds such as the lemniscate, it turns out that there are certain immersed but nonembedded submanifolds for which it holds. To distinguish them, we introduce the following definition:

**Definition 168.** *If $M$ is a smooth manifold and $S \subseteq M$ is an immersed submanifold, then $S$ is said to be **weakly embedded in** $M$ if every smooth map $F\colon N \to M$ whose image lies in $S$ is smooth as a map from $N$ to $S$. (Be aware that weakly embedded submanifolds are called **initial submanifolds** by some authors.)*

Corollary 36 shows that every embedded submanifold is weakly embedded. It follows from Example 66 that the lemniscate is not weakly embedded. However, the dense curve on the torus is weakly embedded (see Problem 5-13 from [Lee, 2013]).

Using the preceding results about restricting maps to submanifolds, we can now prove the promised uniqueness theorem for the smooth manifold structure on an embedded submanifold:

**Theorem 118.** *Suppose M is a smooth manifold and $S \subseteq M$ is an embedded submanifold. The subspace topology on S and the smooth structure described earlier in* Theorem 113 *are the only topology and smooth structure with respect to which S is an embedded or immersed submanifold.*

*Proof.* See proof on [Lee, 2013, p. 114]. □

Thanks to this uniqueness result, we now know that a subset $S \subseteq M$ is an embedded submanifold if and only if it satisfies the local slice condition, and if so, its topology and smooth structure are uniquely determined. Because the local slice condition is a local condition, if every point $p \in S$ has a neighborhood $U \subseteq M$ such that $U \cap S$ is an embedded $k$-submanifold of $U$, then $S$ is an embedded $k$-submanifold of $M$.

Theorem 118 is false in general if $S$ is merely immersed; but we do have the following uniqueness theorem for the smooth structure of an immersed submanifold once the topology is known:

**Theorem 119.** *Suppose M is a smooth manifold and $S \subseteq M$ is an immersed submanifold. For the given topology on S, there is only one smooth structure making S into an immersed submanifold.*

It is certainly possible for a given subset of $M$ to have more than one topology making it into an immersed submanifold (Exercise!). However, for weakly embedded submanifolds we have a stronger uniqueness result:

**Theorem 120.** *If M is a smooth manifold and $S \subseteq M$ is a weakly embedded submanifold, then S has only one topology and smooth structure with respect to which it is an immersed submanifold.*

**Lemma 35 (Extension Lemma for Functions on Submanifolds).** *Suppose M is a smooth manifold, $S \subseteq M$ is a smooth submanifold, and $f \in C^{\infty}(S)$.*

 *a) If S is embedded, then there exist a neighborhood U of S in M and a smooth function $\widetilde{f} \in C^{\infty}(U)$ such that $\widetilde{f}\big|_S = f$.*

 *b) If S is properly embedded, then the neighborhood U in part a) can be taken to be all of M.*

### 3.6.5   The Tangent Space to a Submanifold

Let $M$ be a smooth manifold (with or without boundary), and let $S \subseteq M$ be an immersed or embedded submanifold. Since the inclusion map $\iota \colon S \hookrightarrow M$ is a smooth immersion, at each point $p \in S$ we have an injective linear map $d\iota_p \colon T_pS \to T_pM$. In terms of derivations, this injection works in the following way: for any vector $v \in T_pS$, the image vector $\widetilde{v} = d\iota_p(v) \in T_pM$ acts on smooth functions on $M$ by

$$\widetilde{v}f = d\iota_p(v)f = v(f \circ \iota) = v(f\big|_S).$$



Figure 3.22: The tangent space to an embedded submanifold.

We adopt the convention of identifying $T_pS$ with its image under this map, thereby thinking of $T_pS$ as a certain linear subspace of $T_pM$ (see Figure 3.22). This identification makes sense regardless of whether $S$ is embedded or immersed.

There are several alternative ways of characterizing $T_pS$ as a subspace of $T_pM$. The first one is the most general; it is just a straightforward generalization of a previous proposition:

**Proposition 83.** *Suppose $M$ is a smooth manifold (with or without boundary), $S \subseteq M$ is an embedded or immersed submanifold, and $p \in S$. A vector $v \in T_pM$ is in $T_pS$ if and only if there is a smooth curve $\gamma \colon J \to M$ whose image is contained in $S$, and which is also smooth as a map into $S$, such that $0 \in J$, $\gamma(0) = p$, and $\gamma'(0) = v$.*

The next proposition gives a useful way to characterize $T_pS$ in the embedded case (a nice exercise would be to do Problem $5 - 20$ from [Lee, 2013], which shows that this does not work in the nonembedded case):

**Proposition 84.** *Suppose $M$ is a smooth manifold, $S \subseteq M$ is an embedded submanifold, and $p \in S$. As a subspace of $T_pM$, the tangent space $T_pS$ is characterized by*

$$T_pS = \{v \in T_pM \mid vf = 0 \text{ whenever } f \in C^\infty(M) \text{ and } f|_S = 0\}.$$

If an embedded submanifold is characterized by a defining map, the defining map gives a concise characterization of its tangent space at each point, as the next proposition shows:

**Proposition 85.** *Suppose $M$ is a smooth manifold and $S \subseteq M$ is an embedded submanifold. If $\Phi \colon U \to N$ is any local defining map for $S$, then $T_pS = \ker d\Phi_p \colon T_pM \to T_{\Phi(p)}N$ for each $p \in S \cap U$.*

When the defining function $\Phi$ takes its values in $\mathbb{R}^k$, it is useful to restate the proposition in terms of component functions of $\Phi$.

**Corollary 37.** *Suppose $S \subseteq M$ is a level set of a smooth submersion $\Phi = (\Phi^1, \ldots, \Phi^k) \colon M \to \mathbb{R}^k$. A vector $v \in T_pM$ is tangent to $S$ if and only if $v\,\Phi^1 = \cdots = v\,\Phi^k = 0$.*

# Bibliography

[Carothers, 2000] Carothers, N. (2000). *Real Analysis*. Cambridge University Press, first edition.

[Choquet-Bruhat, 1982] Choquet-Bruhat, Y., D.-M. C. (1982). *Analysis, Manifolds and Physics, Part I: Basics*. North Holland, revised edition.

[Folland, 2007] Folland, G. B. (2007). *Real Analysis: Modern Techniques and Their Applications*. Wiley, second edition.

[Hatcher, 2001] Hatcher, A. (2001). *Algebraic Topology*. Cambridge University Press, first edition.

[Lee, 2011] Lee, J. (2011). *Introduction to Topological Manifolds*. Graduate Texts in Mathematics. Springer, second edition.

[Lee, 2013] Lee, J. (2013). *Introduction to Smooth Manifolds*. Graduate Texts in Mathematics. Springer, second edition.

[Massey, 1991] Massey, W. S. (1991). *A Basic Course In Algebraic Topology*. Graduate Texts in Mathematics. Springer-Verlag New York, LLC, first edition.

[Munkres, 2000] Munkres, J. (2000). *Topology*. Prentice Hall, Inc., second edition.

[Pugh, 2003] Pugh, C. C. (2003). *Real Mathematical Analysis*. Undergraduate Texts in Mathematics. Springer, first edition.

[Rudin, 1964] Rudin, W. (1964). *Principals of Mathematical Analysis*. International Series in Pure and Applied Mathematics. McGraw-Hill, Inc., third edition.

[Spivak, 1971] Spivak, M. (1971). *Calculus on Manifolds: A Modern Approach to Classical Theorems of Advanced Calculus*. Mathematics Mongraphs Series. Westview Press, first edition.

[Stein, 2005] Stein, E., S. R. (2005). *Real Analysis: Measure Theory, Integration, and Hilbert Spaces*. Princeton Lectures in Analysis. Princeton University Press, first edition.

[Stillwell, 2010] Stillwell, J. (2010). *Mathematics and its History*. Undergraduate Texts in Mathematics. Springer, third edition.

# Index