# Mathematical writing
## An undergraduate course

*Franco Vivaldi*
*School of Mathematical Sciences*

Last updated: December 19, 2011

# Preface

This is the textbook for the course *Mathematical Writing,* offered to second year undergraduate students in mathematics at Queen Mary, University of London.

Instructions on writing mathematics are normally given to postgraduate students, to prepare them to write research articles and a thesis. However, there are compelling reasons for providing similar training at undergraduate level.

Every mathematician knows that writing a research paper, or presenting a mathematical result in a seminar, are necessary tests of one's understanding of the material. If a sketched argument has flaws, these flaws will surface as soon as one attempts to convince someone else that the argument is correct.

For the same reason, undergraduate students should be asked to elucidate their thinking in writing. Their first submissions tend to be cryptic collections of symbols, where the onus of decoding the symbols is implicitly left to the teacher or the assessor. All students' written output should evolve rapidly to form coherent —if concise— documents.

Having to use the highly specialised language of mathematics with precision and fluency is a rigorous exercise, which encourages attention to structure and economy of thought. The students' achievements in this area will enrich their learning experience, and provide them with a skill whose importance cannot be overestimated.

The Mathematical writing course has several aims:

- teaching the students how to think rigorously;

- raising the profile of writing in a mathematics degree —it's a vital skill and many students are not good at it (see [13]);

- helping the students consolidate, and perform a synthesis of, the material introduced in the first year;

- teaching some elements of higher mathematics —logic, in particular— using language as a tool;

- training the students in developing and presenting mathematical arguments, with appropriate notation and structure;

- preparing the students for writing a thesis.

The course's writing material is taken mostly from standard first-year university mathematics. This may seem elementary, yet students commented on the 'unexpected depth' required in their thinking, once forced to offer verbal explanations. There is a lot of advanced mathematics hidden within elementary mathematics, which will challenge even the best students.

The development of writing techniques will proceed from the particular to the general, from the small to the large: words, phrases, sentences, paragraphs, to end with short compositions. These may be the introduction of a concept, the proof of a theorem, the summary of a section of a book, the first slides of a presentation.

Chapter 1 is a warm-up, listing some dos and don'ts of writing. In chapter 2 we introduce the basic dictionary pertaining sets, functions, sequences, and equations. Higher mathematics is based on these words, and we begin to use them rigorously. The analysis of mathematical reasoning begins in chapter 3, where we develop some constructs of elementary logic (predicate calculus, quantifiers). This material underpins the expansion of the mathematical dictionary in chapter 4, where basic attributes of real functions are introduced: ordering, symmetry, boundedness, continuity. Mathematical arguments and proofs are studied in detail in the second part of the course (chapters 6–9). Some chapters are dedicated explicitly to writing: chapter 1 gives writing tips, chapter 5 is concerned with mathematical notation, and effectiveness and clarity of exposition.

The symbols $[\not{\mathcal{E}}]$ and $[\not{\mathcal{E}}, n]$ appear often in exercises. They indicate that the written material should contain *no mathematical symbols* (apart from numerals), while the integer $n$ specifies the approximate word length of the assignment. (If $[\not{\mathcal{E}}]$ does not appear, mathematical symbols are intended to be used freely.) In an appropriate context, having to express mathematics without symbols is a most useful exercise. It brings about the discipline needed to use symbols effectively, and it's invaluable for learning how to communicate to an audience of non-experts. Consider the following question:

> *On the plane, I have a circle and a point outside it, and I must find the lines through this point which are tangent to the circle. What shall I do?*

The mathematics is elementary; yet answering the question (with 100–150 words, say) requires organisation, conciseness, and a clear understanding of the structure of the problem. Given the instructional nature of the task, we itemise the answer.

*1. Write down the equation of a line passing through the point. This equation will depend on one parameter, the line's slope, which is the quantity to be determined.*

*2. Adjoin this equation to that of the circle, and solve for the points of intersection of the two curves. After a substitution, you'll end up with a quadratic equation in one unknown, whose coefficients depend on the parameter.*

*3. Equate the discriminant of the quadratic equation to zero, to obtain the desired equation —also quadratic— for the slope.*

*4. If your point lies outside, on, or inside the circle, respectively, you will obtain two distinct real solutions, two identical solutions, or two complex solutions for the slope, respectively. Any geometrical configuration involving vertical lines (infinite slope) will require some care.*

The most challenging exercise of this kind is the MICRO ESSAY, where the synthesis of a mathematical topic has to be performed in a couple of paragraphs, without using any symbols at all. This exercise will prepare the students for writing abstracts, a notoriously difficult task. The mathematical literature is full of abstracts that are too long, or too technical, or use unnecessary symbols.

The timeless, concise book *The elements of style,* by W. Strunk Jr and E. B. White [9] is an ideal complement to the present textbook. Anyone interested in writing should study it carefully.

The available literature on mathematical writing is almost entirely targeted to post-graduate students and researchers. An exception is *How to think like a mathematician*, by K. Houston [6], which is written for beginning university students, and devotes two initial chapters to mathematical writing. The advanced texts include *Mathematical writing,* by D. E. Knuth, T. L. Larrabee, and P. M. Roberts [8], *Handbook of writing for the mathematical sciences* by N. Higham [5], and *A primer of mathematical writing* by S. G. Krantz [7]. Equally valuable is the concise classic text *Writing mathematics well,* by L. Gillman [2]. Unfortunately, this 50-page booklet is out of print, and used copies may command high prices.

This book was inspired by the lecture notes of a course in Logic given by Wilfrid Hodges at Queen Mary in 2005-2006. This course used writing as an essential tool. Wilfrid has been an ideal companion during our decade-long effort to bring writing

Franco Vivaldi
London, 2011.

# Contents

# Chapter 1

# Some writing tips

Before approaching mathematical writing systematically, we consider some general guidelines, illustrated by examples of good and bad practice. Our recommendations are often little more than common sense, or deal with frequent mistakes. Some elementary points concerning style and the use of symbols are raised in the last two sections; they will be developed further in later chapters.

Writing is difficult. The students should return to this chapter repeatedly, to monitor the assimilation of good practice.

## 1.1   Preparation and structure

1. Begin by writing your document in draft form, or at least write down a list of key points. Few people are able produce good writing at the first attempt. Once I marked some short examination essays, and I noticed that half of the students didn't do any preparation. Invariably, their writing was poor.

2. Consider the background of your readers; are they familiar with the meaning of the words you use? It's easy to write a mathematics text that's too difficult; it's almost impossible to write one that's too easy.

3. Form each sentence in your head before writing it down. Then read carefully what you have written. Read it out loud: how does it sound? Have you written what you intended to write? Is it clear? Don't hesitate to rewrite.

4. Split the text into paragraphs. Each paragraph should be about one 'idea', and it should be clear how you are moving from one idea to the next. Be prepared

to re-arrange paragraphs. The first idea you thought of may not have been the best one; the sequence of arguments you have chosen may not be optimal.

5. When you finish writing, consider the opening and closing sentences of your document. The former should motivate the readers to keep reading, the latter should mark a resting place, like the final bars in a piece of music.

6. Word processing has changed the way we write, and often a document is the end-product of several successive approximations. After prolonged editing, one stops seeing things. If you have time, leave your document to rest for a day or two, and then read it again.

## 1.2   Grammar

1. If you are unfamiliar with the basic terminology of grammar (adjective, adverb, noun, pronoun, verb, etc.), look it up in a book, e.g., [1] [9, pp. 89–95].

2. Write in complete sentences. Every sentence should begin with a capital letter, end with a full stop, and contain a subject and a verb. The expression 'A cubic polynomial' is not a sentence, because it doesn't have a verb. It would be appropriate as a caption, or a title, but you can't insert it as it is in the middle of a paragraph.

3. Make sure that the nouns match the verbs grammatically

   BAD:  The set of primes are infinite.

 GOOD:  The set of primes is infinite.

(The verb refers to 'the set', which is singular.)

Make a pronoun agree with its antecedent

   BAD:  Each function should be greater than their derivative.

 GOOD:  Each function should be greater than its derivative.

(The pronoun 'its' refers to 'function', which is singular.)

Do not split infinitives

---

[1] 'abbreviation for the Latin *exempli gratia est,* which means 'for example'.

BAD: We have to again eliminate a variable.

GOOD: We have to eliminate a variable again.

(The infinitive is 'to eliminate'.)

4. Check the spelling: no point in crafting a document carefully, if you then spoil it with spelling mistakes. If you use a word processor, take advantage of a spell checker. These are some frequently misspelled words:

BAD: auxillary, catagory, consistant, correspondance, impliment, indispensible, ocurrence, preceeding, refering, seperate.

These are misspelled mathematical words that I found in mathematics examination papers:

BAD: arithmatic, arithmatric, divisable, infinaty, matrics, orthoganal, orthoginal, othogonal, reciprical, scalor, theorom.

5. Be careful about distinctions in meaning.
Do not confuse *it's* (abbreviation for *it is*) with *its* (possessive pronoun).

BAD: Its an equilateral triangle: it's sides all have the same length.

GOOD: It's an equilateral triangle: its sides all have the same length.

Do not confuse the noun *principle* (general law, primary element) with the adjective *principal* (main, first in rank of importance).

BAD: the principal of induction

BAD: the principle branch of the logarithm

Do not use *less* (of smaller amount, quantity) when you should be using *fewer* (not as many as).

BAD: There are less primes between 100 and 200 than between 1 and 100.

6. Do not use *where* inappropriately. As a relative adverb, *where* stands for *in which* or *to which*; it does not stand for *of which*.

BAD: We consider the logarithmic function, where the derivative is positive.

GOOD: We consider the logarithmic function, whose derivative is positive.

The adverb *when* is subject to similar misuse.

BAD:  A prime number is when there are no proper divisors.

GOOD:  A prime number is an integer with no proper divisors.

7. Do not say *which* when *that* sounds better.  Experiment to decide which is better, and if you can substitute *that* for *which*, do it.[2]  The general rule is to use *which* only when it is preceded by a comma or by a preposition, or when it is used interrogatively. In some cases both pronouns are correct, but have different meaning. *That* is the defining pronoun —it is used to identify an object uniquely— while *which* is non-defining —it adds information to an object already identified.

> The argument that was used above is based on induction.
> [*Specifies which argument.*]
> The following argument, which will be used in subsequent proofs,
> is based on induction.
> [*Adds a fact about the argument in question.*]

8. In presence of parentheses, the punctuation follows strict rules. The punctuation outside parentheses should be correct if the statement in parentheses is removed; the punctuation within parentheses should be correct independently of the outside.

BAD:  This is bad. (Superficially, it looks good).

GOOD:  This is good. (Superficially, it looks the same as the BAD one.)

BAD:  This is bad, (on two accounts.)

GOOD:  This is good (as you would expect).

## 1.3   Style

1. Give priority to clarity over style. Avoid long and involved sentences; break long sentences into shorter ones.

BAD:  We note the fact that the polynomial $2x^2 - x - 1$ has the coefficient of the $x^2$ term positive.

---

[2]American and English writers may have different views on this point.

GOOD: The polynomial $2x^2 - x - 1$ has positive leading coefficient.

BAD: The inverse of the matrix $A$ requires the determinant of $A$ to be non-zero in order to exist, but the matrix $A$ has zero determinant, and so its inverse does not exist.

GOOD: The matrix $A$ has zero determinant, and hence it has no inverse.

2. Prefer the active to the passive voice.

BAD: The convergence of the above series will now be established.

GOOD: We establish the convergence of the above series.

3. Vary the choice of words to avoid monotony. (Use a thesaurus: there is one at `http://thesaurus.reference.com`).

BAD: The function defined above is a function of both $x$ and $y$.

GOOD: The function defined above depends on both $x$ and $y$.

4. Do not use unfamiliar words unless you know their exact meaning.

BAD: A simplistic argument shows that our polynomial is irreducible.

GOOD: A simple argument shows that our polynomial is irreducible.

5. Do not use vague, general statements, to lend credibility to your writing. Avoid emphatic statements.

BAD: Differential equations are extremely important in modern mathematics.

BAD: The proof is very easy, as it makes a quite elementary use of the triangle inequality.

GOOD: The proof uses the triangle inequality.

6. Do not use jargon, or text messages abbreviations: it looks immature rather than 'cool'.

BAD: Spse U subs $x$ into T eq. Wot R T soltns?

7. Enclose side remarks within commas, which is very effective, or parentheses (it gets out of the way). To isolate a phrase, use hyphenation —it really sticks out— or, if you have a word processor, *change font* (**but** don't <u>overdo</u> *it*).

8. Take punctuation seriously. To improve it, begin with [11].

## 1.4   Numbers and symbols

Learning how to combine numbers, symbols and words is one of the aims of this course. We look at some basic conventions, which will be developed further in chapter 5.

1. A sentence containing numbers and symbols must still be a correct English sentence, including punctuation.

   BAD:  $a < b \ a \neq 0$

   GOOD:  Let $a < b$, with $a \neq 0$.

   GOOD:  We find that $a < b$ and $a \neq 0$.

   BAD:  $x^2 - 11^2 = 0.\ x = \pm 11$.

   GOOD:  Let $x^2 - 11^2 = 0$; then $x = \pm 11$.

   GOOD:  The equation $x^2 - 11^2 = 0$ has two solutions: $x = \pm 11$.

2. Omit unnecessary symbols.

   BAD:  Every differentiable real function $f$ is continuous.

   GOOD:  Every differentiable real function is continuous.

3. If you use using small numbers for counting, write them out in full; if you refer to specific numbers, use numerals.

   BAD:  The equation has 4 solutions.

   GOOD:  The equation has four solutions.

   GOOD:  The equation has 127 solutions.

   BAD:  Both three and five are prime numbers.

   GOOD:  Both 3 and 5 are prime numbers.

4. If at all possible, do not begin a sentence with a numeral or a symbol.

   BAD:  $\rho$ is a rational number with odd denominator.

   GOOD:  The rational number $\rho$ has odd denominator.

5. Do not combine operators ($+$, $\neq$, $\Rightarrow$, etc.) with words.

   BAD:  The number $\sqrt{2} - 3/2$ is $< 0$

GOOD: The number $\sqrt{2} - 3/2$ is negative.

BAD: If $a$ is an integer $\Rightarrow a$ is a rational number.

GOOD: If $a$ is an integer, then $a$ is a rational number.

6. Within a sentence, adjacent formulae or symbols must be separated by words.

BAD: Consider $A_n, n < 5$.

GOOD: Consider $A_n$, where $n < 5$.

BAD: Add $p$ $k$ times to $c$.

BAD: Add $p$ to $c$ $k$ times.

GOOD: Add $p$ to $c$, repeating this process $k$ times.

For displayed equations the rules are a bit different, because the spacing between symbols becomes a syntactic element. Thus an expression of the type

$$A_n = B_n, \quad n < 5$$

is quite acceptable.

## Exercises

**Exercise 1.1.** Improve the writing, following the guidelines given in this chapter.

1. $a$ is positive.

2. $X$ is a finite set.

3. We minus the equation.

4. $x^2 + 1$ has no real solution.

5. Two is the only even prime.

6. Suppose $t \neq 0$.

7. When you times it by negative $x$, $<$ becomes $>$.

8. $\sin(\pi x) = 0 \Rightarrow x$ is integer.

9. We have less solutions than we had before.

10. It follows $x - 1 = y^4$.

11. The product of 2 negatives is positive.

12. The set of vertex of pentagons.

13. The set of solutions are all odd.

14. An ellipse is when major and minor axis are the same.

15. The asyntotes of this hyperbola are othogonal.

16. Plug-in the solution in the other equation.

17. The number of primes less that $100 = 25$.

18. When discriminant is $< 0$, you get complex.

19. Let us device a strategy for a proof.

20. The definate integral is when you don't have integration limits.

21. A quadratic function has a single stationery point.

# Chapter 2

# Essential dictionary

In writing mathematics we use words and symbols to describe facts. We need to *explain* the meaning of words and symbols, and to *state* and *prove* the facts. In this chapter we consider the meaning of words and symbols; later in this course we'll be concerned with facts.

Ideally, we ought to explain the meaning of all the words and symbols that we use. But this is impossible: we would need to explain the words used in the explanation, and so on. Instead, we should only explain a word or symbol if our explanation will make it clearer than it was before. Accordingly, we shall call a word or symbol **primitive**[1] if it's suitable to use without explaining its meaning.

An ordinary English word like 'thousand' is obviously primitive, but for more specialised words we must consider the context. When we communicate to the general public, a mathematical term such as **multiplication** can safely be regarded as primitive. Likewise, there should be no need to explain to a mathematician what an **eigenvalue** is, while a number theorist will be familiar with **conductor.** Then there are extremes of specialisation: only a handful of people of this planet will know the meaning of **Hsia kernel.** Finally, terms such as **exceptional set** mean different things in different contexts. Understanding what constitutes an appropriate set of primitives is essential for effective communication of complex knowledge; getting this right is, of course, difficult.

This chapter introduces almost two hundred mathematical words, highlighted in boldface, and provided with accompanying notation. This is our essential mathematical dictionary, built around some basic building blocks: **set**, **function**, **equation**, **sequence**. As we introduce new words, we'll begin to use them. Even though

---

[1]This terminology is due to Pascal (1623-1662).

the writing in this chapter will be limited to short phrases and sentences, it'll soon become clear that using with confidence these words is a lot harder —and more satisfying— than just knowing what they mean.

## 2.1   Sets

A **set** is a collection of *well-defined*, *unordered*, *distinct* objects. (This is the so-called 'naive definition' of a set, due to Cantor[2].)  These objects are called the **elements** of a set, and a set is determined by its elements. We may write

> *The set of all odd integers*
> *The set of vertices of a pentagon*
> *The set of differentiable real functions.*

In simple cases, a set can be defined by listing its elements, separated by commas, enclosed within curly brackets. The expression

$$\{1,2,3\}$$

denotes the set whose elements are the integers 1, 2 and 3. Two sets are equal if they have the same elements:

$$\{1,2,3\} = \{3,2,1\}.$$

(By definition, the order in which the elements of a set are listed is irrelevant.) A **multiset** is a generalisation of a set, whereby its elements need not be distinct: $\{2,1,3,1,3\}$.    Multisets are much less common than sets. Unless this term is mentioned explicitly, one usually adopts the convention that repeated elements are to be ignored:  $\{2,1,3,1,3\} = \{2,1,3\}$.  This convention, implemented in some computer algebra systems, simplifies the definition of sets.

The set $\{\}$ with no elements is called the **empty set**, denoted by the symbol $\emptyset$. The empty set is distinct from 'nothing', it is more like an empty container. For example, the statements

> *This equation has no solutions*
> *The solution set of this equation is empty*

---

[2]Georg Cantor, German mathematician (1845–1918).

have the same meaning.

To assign a symbol to a mathematical object, we use an **assignment statement** (or **definition**), which has the following syntax

$$A := \{1, 2, 3\}. \tag{2.1}$$

This expression assigns the symbolic name $A$ to the set $\{1, 2, 3\}$, and now we may use the former in place of the latter. The symbol ':=' denotes the **assignment operator**. It reads "*becomes*", or "*is defined to be*", rather than "*is equal to*", to underline the difference between assignment and equality (in computer algebra, the symbols = and **:=** are not interchangeable at all!). So we can't write $\{1, 2, 3\} := A$, because the left operand of an assignment operator must be a symbol or a symbolic expression.

The right-hand side of an assignment statement such as (2.1) is a collection of symbols or words that pick out a unique thing, which is called the *definiens* (Latin for 'thing that defines'). The left-hand side is a symbol that will be used to stand for this unique thing, which is called the *definiendum* ('thing to be defined').[3] The definiendum may also be symbolic expression —see below.

While it's very common to use the equal sign '=' also for an assignment, the specialised notation := improves clarity. There are alternative symbols for the assignment operator, namely,

$$\stackrel{\text{def}}{=} \qquad \stackrel{\triangledown}{=}, \tag{2.2}$$

which will make an even stronger point.

To indicate that $x$ is an element of a set $A$, we write

$$x \in A \qquad \text{"}x \text{ is an element of } A\text{"} \qquad \text{"}x \text{ belongs to } A\text{"}.$$

The symbol $\notin$ is used to negate membership. Thus

$$\{7, 5\} \in \{5, \{5, 7\}\} \qquad 7 \notin \{5, \{5, 7\}\}.$$

A **subset** $B$ of a set $A$ is a set whose elements all belong to $A$. We write

$$B \subset A \qquad \text{"}B \text{ is a subset of } A\text{"} \qquad \text{"}B \text{ is contained in } A\text{"}$$

and we use $\not\subset$ to negate set inclusion. For example

$$\{3, 1\} \subset \{1, 2, 3\} \qquad \emptyset \subset \{1\} \qquad \{2, 3\} \not\subset \{2, \{2, 3\}\}.$$

---

[3]These terms are rather heavy, but they are the only ones in use.

Every set has at least two subsets: itself and the empty set. Sometimes these are referred to as the **trivial** subsets. Every other subset —if any— is called a **proper subset.** Motivated by an analogy with $\leqslant$ and $<$, some authors write $\subseteq$ in place of $\subset$, reserving the latter for proper inclusion: $\mathbb{R} \subseteq \mathbb{R}$, $\mathbb{Q} \subset \mathbb{R}$. Proper inclusion is occasionally expressed with the pedantic notation $\subsetneq$.

The number of elements of a set is called the **cardinality**, which is denoted by the prefix #

$$\#\{7, -1, 0\} = 3 \qquad \#A = n.$$

The absolute value symbol $|\cdot|$ is also used to denote cardinality: $|\{7, -1, 0\}| = 3$. When using it, some common sense is needed to avoid any ambiguity. A set is **finite** if its cardinality is finite, and **infinite** otherwise. To indicate that the set $A$ is finite, without disclosing its cardinality, we write

$$\#A < \infty. \tag{2.3}$$

(Characterising the cardinality of infinite sets requires an approach more sophisticated than mere 'counting'.)

Next we consider the words associated to operations between sets. We write $A \cap B$ for the **intersection** of the sets $A$ and $B$: this is the set comprising elements that belong to both $A$ and $B$. If $A \cap B = \emptyset$, we say that $A$ and $B$ are **disjoint,** or have **empty intersection.** The sets $A_1, A_2, \ldots$ are **pairwise disjoint** if $A_i \cap A_j = \emptyset$ whenever $i \neq j$.

We write $A \cup B$ for the **union** of $A$ and $B$, which is the set comprising elements that belong to $A$ or to $B$ (or to both $A$ and $B$).

We write $A \smallsetminus B$ for the **(set) difference** of $A$ and $B$, which is the collection of the elements of $A$ that do not belong to $B$. The **symmetric difference** of $A$ and $B$, denoted by $A \triangle B$, is defined as

$$A \triangle B \overset{\text{def}}{=} (A \smallsetminus B) \cup (B \smallsetminus A).$$

The assignment operator '$\overset{\text{def}}{=}$' (cf. (2.2)) makes it clear that this is a definition. This notation establishes the meaning of $A \triangle B$, which is a symbolic expression rather than an individual symbol. The following examples illustrate the action of set operators.

$$
\begin{aligned}
\{1,2,3\} \cap \{3,4,5\} &= \{3\} \\
\{1,2,3\} \cup \{3,4,5\} &= \{1,2,3,4,5\} \\
\{1,2,3\} \smallsetminus \{3,4,5\} &= \{1,2\} \\
\{1,2,3\} \triangle \{3,4,5\} &= \{1,2,4,5\}.
\end{aligned}
$$

The above are examples of **binary set operators;** they act on sets and have two **operands.** The identities

$$A \cap B = B \cap A \qquad (A \cap B) \cap C = A \cap (B \cap C)$$

express the **commutative** and **associative** properties of the intersection operator. Union and symmetric difference enjoy the same properties, whereas set difference does not.

Let $A$ be a subset of a set $X$. The **complement** of $A$ (in $X$) is the set $X \smallsetminus A$, denoted by $A'$ or by $A^c$. The complement of a set is defined with respect to an **ambient set** $X$. Reference to the ambient set may be omitted, when it's understood. So we may write

*The composite integers are the complement of the primes*

since it's clear that the ambient set is $\mathbb{Z}$.

An **ordered pair** is an expression of the type $(a, b)$, with $a$ and $b$ arbitrary quantities. Ordered pairs are defined by the property

$$(a, b) = (c, d) \qquad \text{if} \qquad a = c \quad \text{and} \quad b = d. \tag{2.4}$$

The ordered pair $(a, b)$ should not be confused with the set $\{a, b\}$, since for pairs order is essential and repetition is allowed. (Ordered pairs may be defined solely in terms of sets —see exercise 13.) Let $A$ and $B$ be sets. We consider the set of all ordered pairs $(a, b)$, with $a$ in $A$ and $b$ in $B$. This set is called the **cartesian product** of $A$ and $B$, and is written as

$$A \times B.$$

Note that $A$ and $B$ need not be distinct; one may write $A^2$ for $A \times A$, $A^3$ for $A \times A \times A$, etc. Because the cartesian product is **associative**, the product of more than two sets is defined unambiguously. Also note that the explicit presence of the multiplication operator '$\times$' is needed here, because the expression $AB$ has a different meaning (see below).

A **partition** of a set $A$ is a collection of pairwise disjoint non-empty subsets of $A$, whose union is $A$. A partition may be described as a **decomposition** of a set into **classes**. For instance, the set $\{\{2\}, \{1, 3\}\}$ is a partition of $\{1, 2, 3\}$, the even and odd integers form a partition of the integers and the plane may be partitioned into concentric circles.

The **power set** $\mathbf{P}(A)$ of a set $A$ is the set of all subsets of $A$. Thus if $A = \{1, 2, 3\}$, then

$$\mathbf{P}(A) = \{\{\}, \{1\}, \{2\}, \{3\}, \{1, 2\}, \{1, 3\}, \{2, 3\}, \{1, 2, 3\}\}.$$

To construct a subset of $A$, we consider each elements of $A$, and we decide whether to include it, or to leave it out. Because any sequence of choices is allowed, if $A$ has $n$ elements, then $\mathbf{P}(A)$ has $2^n$ elements.

## 2.1.1   Defining sets

Defining a set by listing its elements works only for sets of small cardinality. How do we define large or infinite sets? A simple notational device is to use the *ellipsis* '...', which indicates the deliberate omission of certain elements, the identity of which is made clear by the context. For example, the set $\mathbb{N}$ of **natural numbers** is defined as

$$\mathbb{N} := \{1, 2, 3, \ldots\}.$$

Here the ellipsis denotes the omission of all the integers greater than 3. Some authors regard 0 as a natural number, so the definition

$$\mathbb{N} := \{0, 1, 2, 3, \ldots\}$$

is also found in the literature. Both definitions have merits and drawbacks; mathematicians occasionally argue about it, but this issue will never be resolved. So, when using the symbol $\mathbb{N}$, one may need to clarify which version of this set is employed. The set of **integers,** denoted by $\mathbb{Z}$ (from the German *Zahlen*, meaning numbers), can also be defined using ellipses

$$\mathbb{Z} := \{\ldots, -2, -1, 0, 1, 2, \ldots\} \qquad \text{or} \qquad \mathbb{Z} := \{0, \pm 1, \pm 2, \ldots\}.$$

To define general sets we need more powerful constructs. A **standard definition** of a set is an expression of the type

$$\{x : x \text{ has } \mathscr{P}\} \tag{2.5}$$

where $\mathscr{P}$ is some unambiguous property that things either have or don't have. The colon ':' separates out the object's symbolic name from its defining properties. The vertical bar '|' or the semicolon ';' may be used for the same purpose. This expression identifies the set of all objects $x$ that have property $\mathscr{P}$.

Thus the empty set may be defined symbolically as

$$\emptyset \overset{\text{def}}{=} \{x : x \neq x\}.$$

The property $\mathscr{P}$ is '$x$ is not equal to $x$', which is not satisfied by any $x$. Likewise, the cartesian product $A \times B$ of two sets (see section 2.1) may be specified as

$$\{x : x = (a,b) \text{ for some } a \in A \text{ and } b \in B\}.$$

The rule '$x$ has property $\mathscr{P}$' now reads: '$x$ is of the form $(a,b)$ with $a \in A$ and $b \in B$'. The same set may be defined more concisely as

$$\{(a,b) : a \in A \text{ and } b \in B\}.$$

In this variant of the standard definition, the type of objects being considered (ordered pairs) is specified at the outset.

The set $\mathbb{Q}$ of **rational numbers** —ratios of integers with non-zero denominator— is defined as follows

$$\mathbb{Q} := \{\frac{a}{b} : a \in \mathbb{Z}, \ b \in \mathbb{N}, \ \gcd(a,b) = 1)\}. \tag{2.6}$$

The property $\mathscr{P}$ is phrased in such a way as to avoid repetition of elements. This concrete definition is the so-called **reduced form** of rational numbers. The rational numbers may also be defined abstractly, as infinite sets of equivalent fractions —see section 3.3.2.

One might think that in the expression for a set we could choose any property $\mathscr{P}$. Unfortunately this doesn't work, for the following reason, known as the *Russell-Zermelo paradox* (1901). Let $\mathscr{P}$ be the property of being a set that is not a member of itself. Thus the quantity

$$\{3, \{3, \{3, \{3\}\}\}\}$$

has property $\mathscr{P}$, whereas

$$\{3, \{3, \{3, \{3, \ldots\}\}\}\}$$

does not have it. (In the above expression, the nested parentheses must match, so the notation $\{3, \{3, \{3, \{3, \ldots\}$ is incorrect.) Let now $W$ be the set of all sets with property $\mathscr{P}$. Then a set $S$ is a member of $W$ if $S$ is not a member of $S$. Apply this to $W$: $W$ is a member of $W$ precisely if $W$ is not a member of $W$. Impossible!!

Fortunately, we can define a set in such a way that the definition guarantees the existence of the set. A **Zermelo definition** identifies a set $W$, by describing it as

*The set of members of $X$ that have property $\mathscr{P}$*

where the **ambient set** $X$ is given beforehand, and $\mathscr{P}$ is a property that the members of $X$ either have or do not have. In symbols, this is written as

$$W := \{x \in X : x \text{ has } \mathscr{P}\}. \tag{2.7}$$

For example, the expression

*The set of real numbers strictly between 0 and 1*

is a Zermelo definition: the ambient set is the set of real numbers, and we form our set by choosing from it the elements that have the stated property.

EXAMPLE. Turn symbols into words.

$$\{x \in \mathbb{Z} : x \geqslant 0, 2 \mid x\}$$

BAD: The set of integers that are greater than or equal to zero, and such that 2 divides them. [*Robotic.*]

GOOD: The set of non-negative even integers.

Zermelo definitions work because it's a basic principle of mathematics (the so-called *subset axiom*) that for any set $X$ of objects and any property $\mathscr{P}$, there is exactly one set consisting of the objects that are in $X$ and have property $\mathscr{P}$. Some authors leave out the '$\in X$' if it's easy to work out the ambient set $X$. This is appropriate when we deal with the basic sets of arithmetic ($\mathbb{N}, \mathbb{Z}, \mathbb{Q}$, etc.), or with sets that are heavily used in the context.

The definiens of a Zermelo definition has a variable $x$ in it. We could use another letter, but the style of definition requires a variable. In section 3.2 we shall see that the definiens of a Zermelo definition is just a special type of function, called a **predicate.**

## 2.1.2   Sets of numbers

The 'open face' letters $\mathbb{N}$, $\mathbb{Z}$, $\mathbb{Q}$ were introduced in the previous section to represent the natural numbers, the integers, and the rationals, respectively. Likewise, we denote by $\mathbb{R}$ the set of **real** numbers, while the set of **complex** numbers is denoted by $\mathbb{C}$. A 'concrete' definition of $\mathbb{C}$ could be written as

$$\mathbb{C} \stackrel{\text{def}}{=} \{z : z = x + iy, \ i = \sqrt{-1}, \ x, y \in \mathbb{R}\}.$$

The symbol $i$ is called the **imaginary unit,** while $x$ and $y$ are, respectively, the **real part** and the **imaginary part** of the complex number $z$. The sets $\mathbb{R}$ and $\mathbb{C}$ are represented geometrically as the **real line** and the **complex plane**, or **Argand plane**. A plot of complex numbers in the Argand plane is called an **Argand diagram**. We have the chain of proper inclusions

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

We turn to operations involving numbers. The **sum** and **difference** of two numbers $x$ and $y$ are always written $x + y$ and $x - y$, respectively. By contrast, their **product** may be written in several equivalent ways

$$xy \qquad x \cdot y \qquad x \times y, \tag{2.8}$$

and so may their **quotient**

$$\frac{x}{y} \qquad x/y \qquad x : y.$$

(The notation $x : y$ is used mostly in elementary texts.) Do not confuse the product dot '$\cdot$' with the **decimal point** '.', e.g.,

$$3 \cdot 4 = 12 \qquad 3.4 = \frac{17}{5}.$$

The quantity $-x$ is the **negative** of $x$, while the **reciprocal** of $x$, defined for $x \neq 0$, is written as

$$\frac{1}{x} \qquad \text{or} \qquad x^{-1}.$$

The notation for exponentiation is $x^y$, where $x$ is the **base**, and $y$ the **exponent**. If the exponent is a positive integer, then exponentiation is defined as repeated multiplication [4], which may be written symbolically as follows:

$$x^n \stackrel{\text{def}}{=} \underbrace{x \cdots x}_{n} \qquad n \geq 1.$$

The assignment operator $\stackrel{\text{def}}{=}$ indicates that this is a definition. The use of the underbrace is necessary to specify the number of terms in the product, because all terms

---

[4]Defining exponentiation for a general exponent requires the logarithmic and exponential functions.

are identical.  Also note the use of the **raised ellipsis** '$\cdots$' to represent repeated multiplication (or repeated applications of any operator), to be compared with the ordinary ellipsis '$\ldots$', used for sets and sequences (see section 2.4). Thus

$$\underbrace{x \cdots x}_{4} = x \cdot x \cdot x \cdot x \qquad \underbrace{x, \ldots, x}_{4} = x, x, x, x$$

while the notation $x \ldots x$ is incorrect.

In arithmetic, the symbol '$|$' is used for **divisibility.**

$$3 | x \qquad \text{``3 divides } x\text{''} \qquad \text{``x is a multiple of 3''.}$$

A divisor of an integer $n$, which is not 1 or $n$ is called a **proper divisor**. A **prime** is an integer greater than 1 that has no proper divisors. The acronyms **gcd** and **lcm** are used for **greatest common divisor** and **least common multiple.** (The expression **highest common factor** (hcf) —a variant of gcd which is popular in schools— is seldom used in higher mathematics.) Some authors use $(a, b)$ for $\gcd(a, b)$; this is to be avoided, since this notation is already overloaded. Two integers are **co-prime** (or **relatively prime**) if their greatest common divisor is 1.

We now construct new sets from the sets of numbers introduced above.  An **interval** is a subset of $\mathbb{R}$ of the type

$$[a, b] := \{x \in \mathbb{R} : a \leqslant x \leqslant b\}$$

where $a, b$ are real numbers, with $a < b$. This interval is **closed,** that is, it contains its end points. We also have **open** intervals

$$(a, b) := \{x \in \mathbb{R} : a < x < b\}$$

as well as **half-open** intervals

$$[a, b) \qquad (a, b].$$

The notational clash between an open interval $(a, b) \subset \mathbb{R}$ and an ordered pair $(a, b) \in \mathbb{R}^2$ is unfortunate but unavoidable, since both notations are firmly established. For $a = 0$ and $b = 1$ we have the (open, closed, half-open) **unit interval.** A semi-infinite interval

$$\{x \in \mathbb{R} : a < x\} \qquad \{x \in \mathbb{R} : x \leqslant b\}$$

is called a **ray.**    The rays consisting of all positive real or rational numbers are particularly important, and have a dedicated notation

$$\mathbb{R}^+ := \{x \in \mathbb{R}, \ x > 0\} \qquad \mathbb{Q}^+ := \{x \in \mathbb{Q}, \ x > 0\} \qquad (2.9)$$

whereas $\mathbb{Z}^+$ is just $\mathbb{N}$.

Some authors extend the meaning of interval to include also rays and lines, and use expressions such as

$$(-\infty,\infty) \qquad [a,\infty) \qquad (-\infty,b]. \qquad (2.10)$$

Because infinity does not belong to the set of real numbers, a notation such as $[1,\infty]$ is incorrect.

A variant of (2.9) is used to denote non-zero real and rational)numbers

$$\mathbb{R}^* := \{x \in \mathbb{R}, \ x \neq 0\} \qquad \mathbb{Q}^* := \{x \in \mathbb{Q}, \ x \neq 0\}. \qquad (2.11)$$

This notation, while being common, is not universally accepted, and should be used with some care (see section 5.2).

The set $\mathbb{R}^2$ of all ordered pairs of real numbers is called the **cartesian plane,** which is the cartesian product of the real line with itself. If $(x,y) \in \mathbb{R}^2$, then the first component $x$ is called the **abscissa** and the second component $y$ the **ordinate**.

The set $\mathbb{Q}^2 \subset R^2$, the collection of points of the plane having rational coordinates, is called the set of **rational points** in $\mathbb{R}^2$. The set $[0,1]^2 \subset \mathbb{R}^2$ is called the **unit square.** In $\mathbb{R}^3$ we have the **unit cube** $[0,1]^3$, and, for $n > 3$ we have the **unit hypercube** $[0,1]^n \subset \mathbb{R}^n$. These sets have themselves the form of a cartesian product. Not all subsets of a cartesian product are cartesian products of subsets. An example is given by the following subsets of the cartesian plane, are related to the geometrical figure of the circle:

$$\begin{array}{lll} \{(x,y) \in \mathbb{R}^2 : x^2+y^2 = 1\} & \textbf{unit circle} & \\ \{(x,y) \in \mathbb{R}^2 : x^2+y^2 \leqslant 1\} & \textbf{closed unit disc} & (2.12) \\ \{(x,y) \in \mathbb{R}^2 : x^2+y^2 < 1\} & \textbf{open unit disc.} & \end{array}$$

Thus the closed unit disc is the union of the open unit disc and the unit circle.

Let $X$ and $Y$ be sets of numbers. The **(Minkowski) sum** $X+Y$ and **product** $XY$ (also known as **algebraic sum (product))**, are defined as follows

$$X+Y \stackrel{\text{def}}{=} \{x+y : x \in X, y \in Y\} \qquad XY \stackrel{\text{def}}{=} \{xy : x \in X, y \in Y\}$$

with the stipulation that repeated elements are to be ignored. For example, if $X := \{1,3\}$ and $Y := \{2,4\}$, then

$$X+Y = \{3,5,7\} \qquad XY = \{2,4,6,12\}.$$

The expression 'sum of sets' is always understood as a Minkowski sum. In the case of product, it is advisable to use the full expression to avoid confusion with cartesian product.

If $X = \{x\}$ consists of a single element, then we use the shorthand notation $x+Y$ and $xY$ in place of $\{x\}+Y$ and $\{x\}Y$, respectively. For example

$$\frac{1}{2}+\mathbb{N} = \{\frac{1}{2},\frac{3}{2},\frac{5}{2},\ldots\} \qquad 3\mathbb{Z} = \{\ldots,-6,-3,0,3,6,\ldots\}.$$

This notation is economical and effective; it leads to concise statements such as

$$m\mathbb{Z}+n\mathbb{Z} = \gcd(m,n)\mathbb{Z}.$$

(See exercise 5.) Elementary —but significant— applications of this notation are found in **modular arithmetic**. Let $m$ be a positive integer. We say that two integers $x$ and $y$ are **congruent modulo** $m$ if $m$ divides $x-y$. This relation is denoted by[5]

$$x \equiv y\,(\mathrm{mod}\ m).$$

Thus

$$-3 \equiv 7\,(\mathrm{mod}\ 5) \qquad 1 \not\equiv 12\,(\mathrm{mod}\ 7).$$

The integer $m$ is called the **modulus**. The set of integers congruent to a given integer is called a **congruence class**. One verifies that the congruence class of $x$ modulo $m$ is the infinite set $x+m\mathbb{Z}$ (involving sum and product of sets), which is given explicitly as

$$x+m\mathbb{Z} = \{x,x\pm m,x\pm 2m,x\pm 3m,\ldots\}.$$

For example, the odd integers are the integers congruent to 1 modulo 2, which is the congruence class $1+2\mathbb{Z}$.

The set of congruence classes modulo $m$ is denoted by $\mathbb{Z}/m\mathbb{Z}$. If $m = p$ is a prime number, the notation $\mathbb{F}_p$ (meaning 'the field with $p$ elements') may be used in place of $\mathbb{Z}/p\mathbb{Z}$. The set $\mathbb{Z}/m\mathbb{Z}$ contains $m$ elements:

$$\mathbb{Z}/m\mathbb{Z} = \{m\mathbb{Z}, 1+m\mathbb{Z}, 2+m\mathbb{Z},\ldots,(m-1)+m\mathbb{Z}\}.$$

Variants of this notation are used extensively in algebra, where one defines sum/product of more general sets, such as groups and rings.

---

[5]This notation is due to Carl Friedrich Gauss, German mathematician (1777–1855).

## 2.1.3   Writing about sets

We have developed a vocabulary on sets which is sufficient for our purpose. In this section, we present examples of short phrases and sentences which make use of set terminology. Our priorities will be logical and formal accuracy, and conciseness.
   Each of the following phrases defines a specific set.

1. *The set of ordered pairs of complex numbers.*

2. *The set of rational points on the unit circle.*

3. *The set of prime numbers with fifty decimal digits.*

4. *The set of lines in the cartesian plane, passing through the origin.*

Note that we haven't used any symbol. The set in item 1 is just $\mathbb{C}^2$. In item 2, among the infinitely many points of the unit circle, we are interested in those having rational co-ordinates. There is no difficulty in writing this set symbolically

$$\{(x,y) \in \mathbb{Q}^2 : x^2 + y^2 = 1\}$$

although its properties are not obvious from the definition. We can see that this set is non-empty (the points $(0 \pm 1)$, $(\pm 1, 0)$ belong to it), but is it infinite? This example illustrates the power of a verbal definition. Item 3, which defines a subset of $\mathbb{N}$, makes an even stronger point. This set must be extremely large, but how can one show that it is non-empty? In set 4, each line counts as a single element, rather than an infinite subset of the plane (lest our set of lines would be the whole of $\mathbb{R}^2$). The symbolic definition of this set is awkward (see section 2.3); to simplify it, we'll consider suitable **representations** of this set.

It is possible to specify a *type* of set, without revealing its precise identity. In each of the following sets there is at least one unspecified quantity.

*The set of fractions representing a given rational number.*

*The set of divisors of an odd integer.*

*A proper infinite subset of the unit circle.*

*The sum of two finite sets of real numbers.*

*A finite set of consecutive integers.*

*The set of all partitions of a set.*

*A set of partitions of the natural numbers.*

Next we define sets, first with both words and symbols, and then with words only. One should consider the relative merits of the two formulations.

*Let $X = \{3\}$.*
*The set whose only element is the integer 3.*

*Let $X = \{m\}$, for some integer $m$.*
*A set whose only element is an integer.*

*Let $m \in \mathbb{Z}$, and let $X$ be a set such that $m \in X$.*
*A set which contains a given integer.*

*Let $X$ be a set such that $X \cap \mathbb{Z} \neq \emptyset$*
*A set which contains at least one integer.*

*Let $X$ be a set such that $\#(X \cap \mathbb{Z}) = 1$*
*A set which contains one and only one integer.*

In defining the symbol $X$, the combination of the verb 'let' and the equal sign replace an assignment operator. An expression of the type '*Let $X \stackrel{\triangledown}{=} \{3\}$*' would be overloaded.

The distinction between definite and indefinite articles is essential, the former describing a unique object, the latter a class of objects. In the following phrases, a change in one article, highlighted in boldface, has resulted in a logical mistake.

BAD:  A proper infinite subset of **a** unit circle.

BAD:  **A** set whose only element is the integer 3.

BAD:  **The** set whose only element is an integer.

BAD:  **The** set which contains one and only one integer.

As a final exercise, we express some elementary geometric facts using set terminology.

*The intersection of a line and a conic section has at most two points.*

*The set of rational points in an open interval is infinite.*

*A cylinder is the cartesian product of a segment and a circle.*

*There is no finite partition of a triangle into squares.*

> *The intersection of a nested set of open discs may be empty.*

We encourage the reader to re-visit known mathematics from an advanced standpoint, using set language.

## 2.2  Functions

Functions are everywhere. Every time a process transforms a mathematical object into another object, there is a function in the background. We begin with a definition.

A **function** consists of two sets together with a rule[6] that assigns to *each* element of the first set a *unique* element of the second set. The first set is called the **domain** of the function and the second set is called the **co-domain**.    A function whose domain is a set $A$ may also be called a function **over** $A$ or a function **defined on** $A$. The terms **map** or **mapping** are synonymous with function. The term **operator** is used to describe certain types of functions (see below).

A function is usually denoted by a single letter or symbol, e.g., $f$. If $x$ is an element of the domain of a function $f$, then the **value of** $f$ **at** $x$, denoted by $f(x)$ is the unique element of the co-domain that is assigned to $x$ by the rule defining $f$. The notation

$$f : A \to B \qquad x \mapsto f(x) \tag{2.13}$$

indicates that $f$ is a function with domain $A$ and co-domain $B$ that **maps** $x \in A$ **to** $f(x) \in B$. The symbol $x$ is the **variable** or (the **argument**) of the function. The symbols $\to$ and $\mapsto$ have a different meaning, and should not be confused. The function

$$I_A : A \to A \qquad x \mapsto x$$

is called the **identity** (**function**) on $A$. When explicit reference to the set $A$ is unnecessary, the identity is also denoted by $Id$ or $\mathbf{1}$.

All symbolic names in a function definition are inessential; the two expressions

$$f : \mathbb{R} \smallsetminus \{0\} \to \mathbb{R} \quad x \mapsto \frac{1}{x} \qquad x : \mathbb{R} \smallsetminus \{0\} \to \mathbb{R} \quad f \mapsto \frac{1}{f}$$

define exactly the same function (even though the rightmost expression breaks just about every rule concerning mathematical notation —see section 5.2).

---

[6]Below, we'll replace the term 'rule' with something more rigorous.

Functions of several variables are defined over cartesian products of sets. For example, the function

$$f : \mathbb{Z} \times \mathbb{Z} \to \mathbb{N} \qquad\qquad (x, y) \mapsto \gcd(x, y)$$

depends on two integer arguments, and hence it's defined over the cartesian product of two copies of the integers. This definition requires a value for $\gcd(0,0)$, which normally is taken to be zero.

Let $f : A \to B$ be a function. The set

$$\{(x, f(x)) \in A \times B : x \in A\} \qquad\qquad (2.14)$$

is called the **graph** of $f$. So a function is completely specified by three sets: domain, co-domain and graph. We can now re-write the definition of a function, disposing of the rather vague term 'rule' appearing in our original definition, and replacing it with a graph. We write a formal definition, and use a layout appropriate for it.

> DEFINITION.    *A* **function** $f$ *is a* **triple** $(X, Y, G)$ *of non-empty sets. The sets $X$ and $Y$ are arbitrary, while $G$ is a subset of $X \times Y$ with the property that for every $x \in X$ there is a unique pair $(x, y) \in G$. The quantity $y$ is called the* **value of the function at** $x$, *denoted by $f(x)$.*

We see that, besides sets, the definition of a function requires the constructs of order pair and triple. It turns out that these quantities can be defined solely in terms of sets (see exercise 13). So, to define functions, all we need are sets, after all.

Given a function $f : A \to B$, the set

$$\{f(x) : x \in A\}$$

is called the **image** of $A$ under $f$, denoted by $f(A)$. Clearly, $f(A) \subset B$, and $f(A)$ is the smallest set that can serve as co-domain for $f$. Thus $\sin(\mathbb{R})$ is the closed interval $[-1, 1]$. The image $f(X)$ of any subset $X$ of $A$ is defined in a similar manner. The term **range** is an alternative to **image**. This term is sometimes used to mean **co-domain**, thereby creating a potential ambiguity. A **constant** is function whose image consists of a single point.

Formally, the notation $f(X)$ is inconsistent, because we have stipulated that the argument of a function is an element of the domain, not a subset of it. Indeed, in

computer algebra, the quantities $f(x)$ and $f(X)$ are represented by different constructs, e.g., `f(x)` and `map(f,X)` with Maple. However, this suggestive notation is widely used.

A function is said to be **injective** (or **one-to-one**) if distinct points of the domain map to distinct points of the co-domain. A function is **surjective** (or **onto**) if $f(A) = B$, that is, if the image of the domain coincides with the co-domain. A function that is both injective and surjective is said to be **bijective.**

For any non-empty subset $X$ of the domain $A$, we define the **restriction of $f$ to $X$** as

$$f|_X : X \to B \qquad\qquad x \mapsto f(x).$$

Given two functions $f : A \to B$ and $g : B \to C$, we define the new function

$$g \circ f : A \to C \qquad\qquad x \mapsto g(f(x))$$

called the **composition** of $f$ and $g$. The notation $g \circ f$ reminds us that $f$ acts before $g$. The image of $x$ under $g \circ f$ is denoted by $(g \circ f)(x)$, where the parentheses isolate $g \circ f$ as the function's symbolic name. The notation $g \circ f(x)$ is a dangerous hybrid, and should be avoided.

If $f : A \to B$ is a bijective function, then the **inverse** of $f$, denoted by $f^{-1}$, is the function $f^{-1} : B \to A$ such that

$$f^{-1} \circ f = I_A \qquad\qquad f \circ f^{-1} = I_B$$

where $I_{A,B}$ are the identities in the respective sets. A function is said to be **invertible** if its inverse exists. If $f : A \to B$ is injective, then we can always define its inverse by restricting its domain to $f(A)$, if necessary. Let $f : A \to B$ be a function, and let $C$ be a subset of $f(A)$. The set of points

$$f^{-1}(C) \overset{=}{\nabla} \{x \in A : f(x) \in C\} \tag{2.15}$$

is called the **inverse image** of the set $C$.

Because the definition of inverse image of a set under a function does not involve the inverse function, the inverse image exists even if the inverse function does not exist. Consider the expressions

$$f^{-1}(y) \qquad f^{-1}(\{y\}).$$

If $y$ belongs to the image of $f$, then the second expression is well-defined, while the first is defined only if $f$ is invertible. For instance, the sine function, as a real

function, is not invertible, although it becomes invertible when restricted to the interval $[-\pi/2, \pi/2]$. The following example illustrate the two scenarios.

$$\sin^{-1}(\{1\}) = \frac{\pi}{2} + 2\pi\mathbb{Z} \qquad \arcsin(1) = \frac{\pi}{2}.$$

Things get even more confusing when the reciprocal $f(x)^{-1}$ of $f(x)$ comes into play (which exists as long as $f(x)$ is non-zero). The quantities $f^{-1}(x)$ and $f(x)^{-1}$ are unrelated; for instance, in the appropriate domain, we have

$$\sin^{-1}(x) = \arcsin(x) \qquad\qquad \sin(x)^{-1} = \csc(x).$$

As we did previously for sets, we define some functions with short phrases.

1. *The integer function that squares its argument.*

2. *The function that counts the number of primes smaller than a given natural number.*

3. *The function that returns 1 if its argument is rational, and 0 otherwise.*

4. *The function that gives the distance between two points on the unit circle, measured along the circumference.*

Item 2 is a much-studied function in number theory. We surmise that the function in item 3 is defined over the real numbers. The image of the function in item 4 is the closed interval $[0, \pi]$.

With a judicious use of definite and indefinite articles, we can specify a function's type, without committing ourselves to a specific object.

1. *The inverse of a trigonometric function.*

2. *The composition of a function with itself.*

3. *An integer-valued injective function.*

4. *A function which coincides with its own inverse.*

In item 2, we infer that the function maps its domain into itself. Functions of type 4 are called **involutions.**

We will return to writing about functions in chapter 4.

## 2.3   Representations of sets

To be able to work with an **abstract set**, a concrete **representation** of it is needed. Representing a set consists in identifying its elements with a collection of familiar objects, such as vectors or matrices. This identification gives a description of a set in terms of another set. For instance, a representation provide the data structures needed for computer implementation.

More precisely, two sets $A$ and $B$ are said to be **equivalent** (written $A \sim B$) if there is a one-to-one correspondence between the elements of $A$ and the elements of $B$, namely, if there exists a bijective function $f : A \rightarrow B$. Equivalent sets have the same cardinality, and vice-versa. The cardinality of infinite sets is characterised using this equivalence. A set equivalent to $\mathbb{N}$ is said to be **countable,** or **countably infinite**. For instance, $\mathbb{Z}$ is countable, and so is $m\mathbb{Z}$, for any $m \in \mathbb{N}$.

A **representation** of a set $A$ is any set $B$ which equivalent to $A$. (This is the most general acceptation of the term representation. Often, representations are based on a stronger notion of equivalence than the one given above.)

For instance, the open unit interval and the real line are equivalent, as established by the bijective function

$$f : \mathbb{R} \rightarrow (0,1) \qquad x \mapsto \frac{1}{\pi} \arctan(x) + \frac{1}{2}. \qquad (2.16)$$

Likewise, the exponential function establishes the equivalence $\mathbb{R} \sim \mathbb{R}^{+}$.

We consider some representation problems. As a first example, let $L$ be the set of lines in the plane passing through a given point $(a,b)$. This set is infinite. Each element $\lambda$ of $L$ is an infinite subset set of $\mathbb{R}^2$, which we write symbolically as

$$\lambda = \left\{ (x,y) \in \mathbb{R}^2 : y = b + s(x-a) \right\}$$

where $s$ is a real number representing the line's slope. The line $x = a$ is not of this form, and must be treated separately. Collecting all lines together, we obtain a symbolic description of $L$

$$L = \left\{ \left\{ (x,y) \in \mathbb{R}^2 : y = b + s(x-a) \right\} : s \in R \right\} \cup \left\{ \left\{ (a,y) \in \mathbb{R}^2 : y \in \mathbb{R} \right\} \right\}.$$

The simplicity of the verbal definition seems to have drowned in a sea of symbols.

We look for a set equivalent to $L$, with a more legible structure. An obvious simplification results from representing $L$ as a set of **cartesian equations**

$$L \sim \{ y = b + s(x-a) : s \in \mathbb{R} \} \cup \{ x = a \}.$$

We have merely replaced the solution set of an equation with the equation itself. This identification provides the desired bi-unique correspondence between the two sets. We can simplify further. Because $a$ and $b$ are given, there is no need to specify them explicitly: it suffices to give the (possibly infinite) value of the slope. Alternatively, we could identify a line by an angle between 0 and $\pi$, measured with respect to some reference axis passing through the point $(a,b)$. The essence of our set is now evident:

$$L \sim \mathbb{R} \cup \{\infty\} \sim [0, /\pi).$$

The interval $[0, \pi)$ is half-open because both end-points correspond to the same line. The rightmost equivalence may be achieved with a transformation of the type (2.16), where the included end-point 0 corresponds to the point at infinity.

As a second example, let us consider the set $S$ of open segments in the plane. Each segment is identified by its end-points, and each end-point is specified by a pair of real numbers. It would seem that $S \sim \mathbb{R}^4$, but the correspondence between the two sets is not bi-unique, because interchanging the end-points leads to the same segment. Furthermore, if the end-points are the same, we obtain the empty set, not a segment.

Rather than removing from $\mathbb{R}^4$ the unwanted points, we change representation. We identify a segment via its mid-point (a pair of real numbers), length (a positive real number), and orientation (an angle between 0 and $\pi$). We see that

$$S \sim \tilde{S} \qquad \text{where} \qquad \tilde{S} := \mathbb{R}^2 \times \mathbb{R}^+ \times [0, \pi).$$

Consider now the subset $U$ of $S$ consisting of all segments of unit length. Using our representation, we can write

$$U \sim \{(c, r, \theta) \in \tilde{S} : r = 1\} \sim \mathbb{R}^2 \times [0, \pi).$$

In the last equivalence, we have removed the idle variable $r$, whose value is fixed.

## 2.4 Sequences

A **sequence** is an ordered list of mathematical objects —not necessarily distinct— called the **terms** (or the **elements**) of the sequence. The terms of a sequence are represented by a common symbol, and each terms is identified by an integer **subscript**

$$(a_1, a_2, \ldots, a_n) \qquad \qquad (a_1, a_2, \ldots). \qquad \qquad (2.17)$$

Here the common symbol is $a$, and the integer values assumed by the subscript begin from 1. The quantity $a_1$, reads *"a sub 1"*, etc., and subscripts may also begin from 0, or from anywhere. The expression on the left denotes a finite sequence, the one on the right suggests that the sequence is infinite.

The **length** of a sequence is the number of its elements. Two sequences are the same if they have the same length, and if their corresponding terms are equal. If $k$ is an unspecified integer, then $a_k$ is called the **general term** of the sequence.

For example, the sequence of primes

$$(p_1, p_2, p_3, \ldots) = (2, 3, 5, \ldots)$$

is infinite. The general term $p_k$ is the $k$th prime number.

In addition to (2.17), there are several notations for sequences, which display the general term alongside information about the subscript range:

$$(a_k)_{k=1}^n \qquad (a_k)_1^n \qquad (a_k)_{k=1}^\infty \qquad (a_k)_{k \geqslant 1} \qquad (a_k). \qquad (2.18)$$

There are also the **doubly-infinite sequences,** where the subscript runs through all the integers

$$(a_k)_{k=-\infty}^\infty = (\ldots, a_{-1}, a_0, a_1, \ldots).$$

In section 5.2 we shall discuss the usage of the various notations for sequences.

The ellipsis is much used in sequence notation, but it must be used wisely. Thus the general term of the sequence of monomials

$$(2x, 2x^2, 2x^3, \ldots)$$

is clearly equal to $2x^k$. However, the expression

$$(3, 5, 7, \ldots)$$

is ambiguous, because there are several plausible alternatives for the identity of the omitted terms, such as $(9, 11, 13, \ldots)$ or $(11, 13, 17, \ldots)$. In the former case, we resolve the ambiguity by inserting the general term

$$(3, 5, \ldots, 2k+1, \ldots);$$

In the latter, we need an accompanying sentence.

A **subsequence** of a sequence $(a_k)$ is any sequence that is obtained from $(a_k)$ by deleting terms. For instance, the primes that give remainder 1 upon division by 4 form a subsequence of the sequence of primes

$$(5, 13, 17, 29, \ldots).$$

Some types of sequences have a specialised terminology. We have seen that a two-element sequence may be called an **(ordered) pair**, and a three-element sequence a **triple.** Occasionally one sees the terms **quadruple** or **quintuple** (I wouldn't go much beyond that), while an $n$-element sequence may be called an $n$-**uple**. A finite sequence of numbers may be called a **vector**, in which case we speak of **dimension** rather than length.

As an exercise in notation for sequences, let us define the $n$-dimensional sphere. Its elements are $n$-dimensional vectors, subject to a constraint. The notation employs a combination of ordinary and raised ellipses (compare with the definition (2.12) of the unit circle)

$$\{(x_1,\ldots,x_n) \in \mathbb{R}^n : x_1^2 + \cdots + x_n^2 = 1\} \qquad n\text{-\textbf{dimensional unit sphere}}$$

An infinite sequence $(a_1, a_2, \ldots)$ represents a **function** defined over the natural numbers. If the elements of the sequence belong to a set $A$, then such a function is defined as

$$a : \mathbb{N} \to A \qquad k \mapsto a_k.$$

We see that in the expression $a_k$, the symbol $a$ is the function's name, the subscript $k$ is an element of the domain, and $a_k = a(k) \in A$ is the value of the function at $k$. This interpretation clarifies the meaning of expressions such as $a_{k^2}$: this the composition of two functions, much like $\sin(x^2)$.

Finally, we develop the notation for **sets of sequences**. Let $A$ be a set. For given integer $n$, let us first consider the set of all *finite* sequences of elements of $A$, with $n$ elements. This set is the cartesian product $A^n$ of $n$ copies of the set $A$: the first element $a_1$ of a sequence is chosen from the first copy of $A$, the second element from the second copy of $A$, and so on. For instance, the set $\{0,1\}^n$ is the set of all binary sequences with $n$-digits, while $\mathbb{Q}^n$ is the set of all $n$-uples of rational numbers. By the same token, the expression

$$\bigcup_{n \geq 1} \mathbb{Z}^n$$

provides a concise symbolic characterisation of the set of all finite integer sequences.

It follows from the above considerations that the set of all infinite sequences of elements of $A$ has the structure of a cartesian product with infinitely many terms. It seems natural to denote such a set by $A^\infty$. However, the idiomatic notation $A^{\mathbb{N}}$ is more common, due to its greater flexibility. Thus $A^{\mathbb{Z}}$ denotes the set of doubly-infinite sequences of elements of $A$, and one even finds $A^{\mathbb{Z}^2}$ for the set of sequences with two indices!

## 2.4.1 Some constructions involving sequences

Let $(A_k)$ be sequence of sets, which may be finite or infinite. The binary set operations of union and intersection generalise to an arbitrary number of operands as follows

$$\bigcup_k A_k = A_1 \cup A_2 \cup \cdots.$$

This expression denotes the set of elements belonging to at least one of the sets $A_k$. As with any binary operator, the implied intersections are represented by the **raised ellipsis.** Likewise the expression

$$\bigcap_k A_k = A_1 \cap A_2 \cap \cdots$$

represents the set of elements belonging to every one of the sets $A_k$. A sequence of sets is **descending** (or **nested**) if

$$A_1 \supset A_2 \supset A_3 \supset \cdots$$

and **ascending** if

$$A_1 \subset A_2 \subset A_3 \subset \cdots.$$

We describe the expression

$$\mathbb{Z} \supset 2\mathbb{Z} \supset 2^2\mathbb{Z} \supset \cdots \supset 2^k\mathbb{Z} \supset \cdots.$$

as follows:

> *An infinite descending chain of sets of even integers, whose intersection consists of a single integer.*

(Think about it.)

Given a *finite* sequence of numbers $(a_1, \ldots, a_n)$, we form the sum and the product of its elements

$$\sum_{k=1}^{n} a_k = a_1 + a_2 + \cdots + a_n \qquad \prod_{k=1}^{n} a_k = a_1 \times a_2 \times \cdots \times a_n. \qquad (2.19)$$

Again, the raised ellipsis represent repeated additions and multiplications, respectively. The symbol $\sum$ is called the **summation symbol.** The subscript $k$ is the **index of summation**, while 1 and $n$ are, respectively, the **lower bound** and **upper bound**

of summation. The quantity $a_k$ is the **general term** of the sum. The integer sequence $(1, 2, \ldots, n)$, specifying the values assumed by the index of summation, is called the **range of summation.** The symbol $\prod$ is called the **product symbol,** and all terminology introduced for sums extends with obvious modifications to products.

If the number of summands is infinite, then the sum in (2.19) is called a **series**. If the limit of finite sums

$$\lim_{n \to \infty} \sum_{k=1}^{n} a_k$$

exists, then such a limit is called the **sum of the series,** and the series is said to **converge**. Otherwise the series **diverges.** If a series has non-negative terms, then convergence is sometimes expressed with the suggestive notation (cf. (2.3))

$$\sum_{k \geq 0} a_k < \infty. \tag{2.20}$$

An infinite product is called just that, and its convergence is defined as the limit of a sequence of finite products. Infinite products are often written in the form

$$\prod_{k \geq 0} (1 + a_k)$$

because convergence requires that $a_k \to 0$. Sum and product notation will be considered again in section 5.3.

Let $A$ be a set of numbers, and let $(a_0, \ldots, a_n)$ be a finite sequence of elements of $A$ with $a_n \neq 0$. A **polynomial** over $A$ in the **indeterminate** $x$ is an expression of the type

$$a_0 + a_1 x + a_2 x^2 + \cdots + a_n x^n. \tag{2.21}$$

The elements of the sequence are called the **coefficients** of the polynomial, and the integer $n$ is its **degree.** The coefficients $a_0$ and $a_n$ are called, respectively, the **constant** and the **leading coefficient.** Each addendum in a polynomial is called a **monomial.** A polynomial of degree two is said to be **quadratic**; then we have **cubic**, **quartic**, **quintic**. The set of all polynomials over the set $A$ with indeterminate $x$ is denoted by $A[x]$. For example

$$x^2 - x - 1 \in \mathbb{Z}[x] \qquad\qquad \frac{1}{2} - y^3 \in \mathbb{Q}[y].$$

A **multivariate polynomial** is a polynomial in more than one indeterminate. (The term **univariate** is used to differentiate from multivariate.)

$$x^2 y^2 - \frac{1}{2} x^4 - x y^3 \in \mathbb{Q}[x, y].$$

The **total degree** of each monomial is the sum of the degrees of the indeterminates, and degree of a polynomial is the largest total degree among the monomials with non-zero coefficient. A multivariate polynomial is **homogeneous** if all monomials have the same total degree. The expression above may be described as

> *A homogeneous quartic polynomial in two indeterminates, with rational coefficients.*

If we replace the finite sum (2.21) with an infinite sum, we obtain a **(formal) power series**.

$$\sum_{k\geq 0} a_k x^k \tag{2.22}$$

This is a 'polynomial of infinite degree'. The attribute 'formal' is used if we are not concerned with assigning specific values to the indeterminate $x$. In this case the power series, like a polynomial, is an algebraic object. For the values of $x$ for which the series (2.22) converges, the power series represents a function.

EXAMPLE. Explain what is a polynomial. [∉]

> *A polynomial is a finite sum. Each term, called a monomial, is the product of a coefficient (typically, a real or complex number) and one or more indeterminates, each raised to some non-negative integer power.*

## 2.5   Equations

Let $f$ and $g$ be functions with the same domain $X$ and co-domain $Y$. In the most general setting, nn **equation** (**on** or **over** $X$) is an expression of the type

$$f(x) = g(x). \tag{2.23}$$

The quantity $x$ is the equation's **unknown**. The expression (2.23) defines a property that each point $x \in X$ either has or doesn't have[7]. This leads to the Zermelo definition of a set

$$\{x \in X : f(x) = g(x)\}$$

which is called the **solution set** of the equation $f(x) = g(x)$.

---

[7]In chapter 3 we shall see that an equation is a special type of **predicate**, which is, in turn, a special type of **function**.

For instance, the expression

$$x^2 - 3x + 1 = -(1 + 3x)$$

is an equation. As an equation over $\mathbb{R}$, its solution set is empty, while over $\mathbb{C}$, its solution set is $\{\sqrt{-2}, -\sqrt{-2}\}$. Thus the solution set of an equation depends on the ambient set.

An expression of the type

$$\begin{cases} f_1(x) = g_1(x) \\ f_2(x) = g_1(x) \\ \quad\vdots \\ f_n(x) = g_n(x) \end{cases} \qquad x = (x_1, \ldots, x_m) \tag{2.24}$$

where all functions have the same domain and co-domain, is called a **system of $n$ simultaneous equations in $m$ unknowns.** The solution set of a system of equations is the intersection of the solution sets of the individual equations (see example 3.2, on page 58).

Equations are very general objects, and need not be associated to 'numbers'. For instance, let $A$ be a set, let $a, b$ be distinct elements of $A$, and let $\mathbf{P}(A)$ be the **power set** of $A$ (see section 2.1). We consider the following **set equation**

$$x \cap \{a\} = x \cap \{b\} \qquad a \neq b. \tag{2.25}$$

With the notation introduced above, we have

$$X = Y = \mathbf{P}(A) \qquad f(x) = x \cap \{a\} \qquad g(x) = x \cap \{b\}.$$

The solution set of the equation is readily seen to be the collection of the subsets $x$ of $A$ which do not contain $a$ or $b$.

If the co-domain $Y$ of $f$ and $g$ is a **group** with respect to addition (e.g., $Y$ is a set of numbers, $Y = \mathbb{Z}, \mathbb{Q}, \mathbb{R}, \mathbb{C}$), then, by replacing $f$ by $f - g$, we can write equation (2.23) in the simpler form

$$f(x) = 0. \tag{2.26}$$

Not all equations can be put in this form, for example equation (2.25) (see exercise 12). In an equation of the form (2.26), the zero on the right-hand side is the zero element of the group, which is not necessarily the number $0$ (see below). An element of the solution set of this equation is called a **zero** of $f$, but if $f(x)$ is a polynomial, we speak of a **root** of $f$. If $f(x) = 0$, then we also say that $f$ **vanishes at** $x$. A

function $f$ **vanishes identically** on a set, if it vanishes at every point of this set. For example, the real function $x \mapsto \sin(\pi x)$ vanishes identically on $\mathbb{Z}$.

A **differential equation** is an equation involving derivatives of the unknown

$$\frac{dx}{dt} + x = 0 \qquad\qquad \frac{dx^2}{dt^2} + t\frac{dx}{dt} + x = 0. \qquad\qquad (2.27)$$

Each of the expressions above is of the form (2.26) for a suitable function $f$. It should be clear that in equation (2.26) the unknown is $x$ (not $t$), and $x = x(t)$ is a **differentiable function** over, say, the real numbers $x : \mathbb{R} \to \mathbb{R}$. The domain of $f$ is the set of differentiable real functions, which is denoted by $C^1(\mathbb{R})$, while the co-domain is the set of continuous real functions, denoted by $C^0(\mathbb{R})$. The solution set is the subset of $C^1(\mathbb{R})$ constituted by all functions $x$ for which $f(x) = 0$. Thus the symbol on the right-hand side of equation (2.26) does not represent the number 0, but rather the zero function $t \mapsto 0$. (Think about it.)

More generally, an equation whose unknown is a function is called a **functional equation**

$$\Gamma(x+1) = x\Gamma(x) \qquad\qquad f = f \circ f \circ f. \qquad\qquad (2.28)$$

There are equations whose unknown are vectors, or matrices, or many other things; in this case 0 represents the neutral element of addition in the appropriate space.

An equation whose solution set is equal to the ambient set is called an **identity, or an indeterminate equation.** Any identity typically reduces to the standard form $0 = 0$. This doesn't mean that identities are trivial: rather they are ephemeral quantities, which express the equivalence of two functions, but which disappear if they are simplified. For instance, the identity

$$x^{2^n} - y^{2^n} = (x-y)\prod_{k=0}^{n-1}(x^{2^k} + y^{2^k})$$

gives the full factorisation of the difference of two monomials whose degree is a power of 2, into the product of polynomials with integer coefficients.

Let us consider the following expressions in two real variables.

$$x + y = y + x \qquad\qquad x + y = 1 - y.$$

The first expression is an identity (the solution set is $\mathbb{R}^2$), representing the commutativity of the addition of real numbers; the second is an equation, whose solution set is a line in $\mathbb{R}^2$.

By restricting the ambient set to the solution set, an equation becomes an identity. For example, the expression $\sin(\pi x) = 0$ is an equation over $\mathbb{R}$ and an identity over $\mathbb{Z}$. There are more meaningful examples of this phenomenon. For instance, let us consider the equation $f(x) = x^5 - x = 0$. If the ambient set is $X = \mathbb{C}$, then the factorisation $x^5 - x = x(x-1)(x+1)(x^2+1)$ shows that the solution set is $\{0, \pm 1, \pm\sqrt{-1}\}$. Consider now the ambient set

$$X = \mathbb{Z}/5\mathbb{Z} = \{0+5\mathbb{Z},\, 1+5\mathbb{Z},\, 2+5\mathbb{Z},\, 3+5\mathbb{Z},\, 4+5\mathbb{Z}\}$$

of congruence classes modulo 5.

We evaluate our function $f$ at all points of $\mathbb{Z}/5\mathbb{Z}$, writing $k$ for $k+5\mathbb{Z}$:

$$
\begin{aligned}
f(0) &= 0^5 - 0 = 0 \equiv 0 \,(\mathrm{mod}\ 5)\\
f(1) &= 1^5 - 1 = 0 \equiv 0 \,(\mathrm{mod}\ 5)\\
f(2) &= 2^5 - 2 = 30 \equiv 0 \,(\mathrm{mod}\ 5)\\
f(3) &= 3^5 - 3 = 240 \equiv 0 \,(\mathrm{mod}\ 5)\\
f(4) &= 4^5 - 4 = 1020 \equiv 0 \,(\mathrm{mod}\ 5).
\end{aligned}
$$

we have shown that, over $\mathbb{Z}/5\mathbb{Z}$, the function $f(x) = x^5 - x$ vanishes identically, so that the equation $f(x) = 0$ is an identity!

EXAMPLE. Explain what is an equation, and its solutions. [¢]

BAD: *An equation is when we equate two functions. The solution is when the functions are the same.*

The inappropriate use of the adverb 'when' is easily spotted (see section 1.2), but there is a more serious flaw. The expression 'equating two functions' means that we seek conditions under which the two functions become the same function. That's not what we had in mind. The equal sign in expression (2.23) is a binary operator, and its operands are not functions, but rather values of functions.

GOOD: *An equation is an expression that identifies the value of two functions at a generic point of their common domain. The solutions of an equation are the points at which the two functions assume the same value.*

(The expression 'equating two functions' may be appropriate for **functional equations**, see (2.28).)

## 2.6 Expressions

The generic term **expression** indicates the symbolic encoding of a mathematical object. For instance, the string of symbols '$2+3$' is a valid expression, and so is '$x \mapsto f(x)$', while '$2+\times 3$' is incorrect.

It would seem that any correct expression should have —in principle, at least— a unique **value**, representing some agreed 'fully simplified' form of the expression. For instance, it could be argued that the value of $\sqrt{2187}$ is $27\sqrt{3}$. Such a value would enable us, among other things, to recognise when two expressions represent the same thing.

This is not so simple. For example, the two expressions

$$\sqrt{3} - \sqrt{2} \qquad \sqrt{5 - 2\sqrt{6}}$$

have the same value, yet there is no compelling reason for choosing one over the other, and our choice will depend on the context. The following well-known identity makes an even stronger point:

$$\frac{1 - x^n}{1 - x} = 1 + x + x^2 + \cdots + x^{n-1}.$$

The right-hand side —the sum of $n$ monomials— is the 'fully simplified' version, while the 'unsimplified' left-hand side involves only four terms.

Given that defining *the* value of expressions proves difficult, we shift our attention to their properties. The most general property is the **type** of value (set, number, function, etc.), which assigns the expression to a certain class. Again, we must exercise some judgement. The expressions

$$1 + 1 \qquad \int_0^\pi \sin(x)dx$$

have the same value, but their structure is so different, that the coincidence of their values seems secondary. Whereas the expression on the left is unquestionably "*a number*", or "*a positive integer*", that on the right is "*a definite integral.*" On the other hand, there may be circumstances in which the reductionist description of the integral as a number is appropriate, for instance when discussing integrability of functions.

The rest of this chapter is devoted to the description of expressions.

## 2.6.1   Levels of description

We develop the idea of successive refinements in the description of a mathematical expression, from the general to the particular. The appropriate level of details to be included will change, depending on the situation. We treat in parallel verbal and symbolic descriptions, as far as it is reasonable to do so.

Let us consider the definition of a set. The coarsest level of description is

$$\{\dots\} \qquad \textit{A set}$$

where the object's type is identified by the curly brackets. The use of the indeterminate article —'a' set rather than 'the' set— reflects our incomplete knowledge.

The next level in specialisation identifies the ambient set

$$\{(x,y) \in \mathbb{Z}^2 : \dots\} \qquad \textit{A set of integer pairs.}$$

Now we begin to build the defining properties of our set

$$\{(x,y) \in \mathbb{Z}^2 : \gcd(x,y) = 1, \dots\} \qquad \textit{A set of pairs of co-prime integers.}$$

The final step gives us complete knowledge

$$\{(x,y) \in \mathbb{Z}^2 : \gcd(x,y) = 1,\ 2|xy\} \qquad \textit{The set of pairs of co-prime integers, with at least one even component.}$$

Accordingly, the indefinite article has been replaced by the definite article. Now both words and symbols describe one and the same object, and one should consider the relative merits of the two presentations. A robotic translation of symbols into words

*The set of elements of the cartesian product of the integers with themselves, whose components have greatest common divisor is equal to 1, and such that 2 divides the product of the components*

while being correct, lacks the synthesis that comes with understanding.

Expression may be **nested,** like boxes within boxes. Let us begin with the expressions

$$(\dots)^2 \qquad \textit{A square}$$
$$\sum \dots \qquad \textit{A sum}$$

We only see the outer structure of these objects. We compose them in two different ways

$$\left( \sum \cdots \right)^2 \qquad \textit{The square of a sum}$$
$$\sum \left( \cdots \right)^2 \qquad \textit{A sum of squares}$$

The first term in each expressions identifies the object's outer layer. There is still one indefinite article in each expression, reflecting a degree of generality. We specialise further:

$$\sum_{n=1}^{\infty} \left( \frac{1}{n} \right)^2 \qquad \textit{The sum of the square of the reciprocal of the natural numbers.}$$

Words or symbols now define a unique object, with three levels of nesting. By contrast, in the nested expression

$$\left( \sum_{n=1}^{\infty} a_n \right)^2 \qquad a_n \in \mathbb{Q} \qquad \textit{The square of the sum of the elements of a rational sequence,}$$

the innermost object —a rational sequence— is still generic.

These examples may give the impression that words and symbols are interchangeable. Not so. Some concepts are best expressed with words, others with symbols, while most situations require a careful combination of the two. For instance, the symbolic expression

$$(1 - x, 2 + x^2, \ldots, n + (-x)^n, \ldots) \tag{2.29}$$

defines a sequence succinctly and unambiguously. Using words, we could refine its description as follows

*A sequence*

*An infinite sequence*

*An infinite sequence of polynomials.*

Increasing further the accuracy of the verbal description seems pointless, since the symbolic expression (2.29) is clearly superior in delivering exact information. On the other hand, with words we can place this expression *in a context,* which is something symbols can't do. We supplement the description given above with additional information, so as to to emphasise different properties.

*An infinite sequence of polynomials*

> *with integer coefficients*
> *with unbounded coefficients*
> *with increasing degree*
> *whose leading term alternates in sign.*

## 2.6.2   Characterising expressions

We expand our dictionary with a list of generic terms of common use, and illustrate their usage with short phrases.

An expression involving numbers, the four arithmetical operations, and raising to an integer or fractional power (extraction of roots), is called an **arithmetical expression.** The value of an arithmetical expression is a number. A combination of rational numbers and square roots of rational numbers is called a **quadratic irrational** or a **quadratic surd**. The following expressions are arithmetical.

$$191861^2 - 3 \cdot 110771^2 = -2$$

*An arithmetic identity, with a surprising cancellation.*

$$\frac{3 + 2\sqrt{2}}{8 - 3\sqrt{7}}$$

*The ratio of two quadratic surds having distinct radicands.*

If indeterminates are present, we speak of an **algebraic expression**

$$\frac{\sqrt[6]{(b - (1/b))^2}}{\sqrt[3]{ab^2 + ab}}$$

*An algebraic expression in two indeterminates.*

**Polynomials** are algebraic expressions (see section 2.4.1) and so are the **rational functions**, namely the ratio of two polynomials. Because polynomials and rational functions do not involve fractional powers of the indeterminates, they may be characterised as **rational expressions**.

$$\frac{1}{x+1} + \frac{1}{x^2+1} + \cdots + \frac{1}{x^n+1}$$

*The sum of finitely many rational functions, with increasing degree.*

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

*A rational expression, obtained by composing a polynomial function with itself several times.*

The following mathematical Russian doll

$$\sqrt{x + \sqrt{x^2 + \sqrt{x^4 + \cdots + \sqrt{x^{2^n}}}}} \qquad n \in \mathbb{N}$$

could be described as

> *A family of algebraic expressions in one indeterminate, involving an arbitrarily large number of nested square roots.*

The functions sine, cosine, tangent, secant, etc., are called **trigonometric functions** (or **circular functions**). A **trigonometric expression** is an expression containing trigonometric functions

$$8\cos(z)^4 - 8\cos(z)^2 + 1 \qquad \textit{A quartic trigonometric polynomial.}$$

Trigonometric functions belong to the larger class of **transcendental functions**, namely functions not definable by an algebraic expression. The exponential and the logarithm are transcendental functions.

An equation defined by algebraic expressions is called an **algebraic equation.** Likewise, we speak of **trigonometric** and **transcendental** equations.

$$x^n - x - 1 = 0 \qquad \textit{An algebraic equation.}$$
$$\cos(x) = \sin(x) \qquad \textit{A trigonometric equation.}$$
$$\log(1 + x) = -x \qquad \textit{A transcendental equation.}$$

The term **analytical expression** is appropriate in the presence of infinite processes:

$$\sqrt{1 + x} = 1 + \frac{1}{3}x - \frac{1}{9}x^2 + \frac{5}{81}x^3 + \cdots \qquad \text{The first few terms of the series expansion of an algebraic function}$$

$$\lim_{n \to \infty} \left(1 + \frac{1}{n}\right)^n = e \qquad \text{Napier's constant as the limit of a sequence of rational numbers}$$

$$2 \prod_{k=1}^{\infty} \frac{(2k)^2}{(2k)^2 - 1} = \pi \qquad \text{An infinite product formula for Archimedes constant.}$$

An **integral expression** is an expression involving integrals

$$\ln(x) = \int_0^x \frac{1}{t}\,dt \qquad \text{An integral expression for the natural logarithm.}$$

The term **combinatorial** is appropriate to expressions involving counting functions, such as the **factorial** function, or the **binomial coefficient**.

$$\frac{1}{2^{2k}}\frac{(2k)!}{(k!)^2}$$  *A rational combination of exponentials and factorials.*

$$\sum_{k=0}^{n}\binom{n+k}{k}$$  *A finite sum of binomial coefficients.*

In chapter 3 we shall deal with the expressions found in logic: the **boolean** expressions.

## Exercises

**Exercise 2.1.** Consider the following topics

  *i*) Prime numbers;      *ii*) fractions;      *iii*) complex numbers.

For each topic **[∉]**

1. Write five short sentences. Each sentence should give a definition or state a fact.

2. Ask five questions. They should have mathematical significance, and preferably possess a certain degree of generality.

    BAD: Is 39 a prime number?  [*Specific and insignificant.*]

   GOOD: Why is 1 not a prime number?

**Exercise 2.2.** Define five interesting finite sets. **[∉]**

BAD: The set of natural number less than 10.

GOOD: The power set of a finite set.

**Exercise 2.3.** The following expressions define sets. Turn words into symbols, using standard or Zermelo definitions.
[*Represent geometrical objects, e.g., planar curves, by their cartesian equations.*]

1. The set of negative odd integers.

2. The complement of the open unit disc in the complex plane.

3. The set of all rational numbers which are not the ratio of two consecutive integers.

4. The set of vectors of unit length in three-dimensional euclidean space.

5. The set of circles in the plane, passing through the origin.

6. The set of hyperbolae in the cartesian plane, whose asymptotes are the cartesian coordinate axes.

7. The set of lines tangent to the unit circle. [*Think about it.*]

**Exercise 2.4.** The following expressions define sets. Turn symbols into words. [∉]

1. $\{x \in \mathbb{Q} : 0 < x < 1\}$

2. $\{1/2n : n \in \mathbb{N}\}$

3. $\{x + y\sqrt{-1} : x, y \in \mathbb{Z}\}$

4. $\{x \in \mathbb{Q} \smallsetminus \mathbb{Z} : 1/x \in \mathbb{Z}\}$

5. $\{x \in \mathbb{Q} \smallsetminus \mathbb{Z} : x^2 \in \mathbb{Z}\}$

6. $\{x \in \mathbb{C} \smallsetminus \mathbb{R} : x^2 \in \mathbb{R}\}$

7. $\{(x, y) \in \mathbb{Z}^2 : \gcd(x, y) > 1\}$

8. $\{(x_1, x_2, x_3) \in \mathbb{R}^3 : x_2 = 0\}$

9. $\{x \in \mathbb{R} : \sin(2\pi x) = 0\}$

10. $\{(x, y) \in \mathbb{R}^2 : \sin(\pi x) = 0\}$.

**Exercise 2.5.** For each item, provide two levels of description: [⨏]

$(i)$  a coarse description, which only identifies the object's type (set, function, polynomial, etc);
$(ii)$  a finer description, which defines the object in question, or characterises its structure.

1. $3^3 + 4^3 + 5^3 = 6^3$

2. $\frac{7}{5} < \sqrt{2} < \frac{17}{12}$

3. $x^3 - x - 1$

4. $(x+y)^3 = x^3 + 3x^2y + 3xy^2 + y^3$

5. $xy > 0$

6. $y - x^2 - x = 0$

7. $\ddot{x} - 3\dot{x} - 2 = 0$

8. $\sin(x - y) = \sin(x)\cos(y) - \cos(x)\sin(y)$

**Exercise 2.6.** Same as in previous problem.

1. $A \cup B = B \cup A$

2. $f(A) \cap f(B)$

3. $f(A \cup B)$

4. $2\mathbb{Z} \smallsetminus 4\mathbb{Z}$

5. $(\mathbb{R} \smallsetminus \mathbb{Q})^2$

6. $(1, 5, 9, \ldots, 4n + 1, \ldots)$

7. $(1, 2^{2k}, 3^{2k}, \ldots, n^{2k})$

8. $\sum_{i=1}^{\infty} |a_i - b_i|, \quad a_i, b_i \in \mathbb{C}$

**Exercise 2.7.** Same as in previous problem.

1. $f : \mathbb{R} \to \mathbb{R}, \quad x \mapsto x + 1$

2. $\int_0^\infty \frac{f(x)}{g(x)} dx$

3. $\frac{d}{dx} f(x)g(x) = \frac{df(x)}{dx} g(x) + \frac{dg(x)}{dx} f(x)$

4. $f(x) = \int_0^x \frac{df(y)}{dy} dy$

5. $\cos(x) = \sum_{k=0}^{\infty} \frac{(-1)^k x^{2k}}{(2k)!}$

6. $\int_0^1 \int_0^1 F(x,y) dx dy$

7. $\frac{\partial F(x,y,z)}{\partial x} + \frac{\partial F(x,y,z)}{\partial y} + \frac{\partial F(x,y,z)}{\partial z}.$

**Exercise 2.8.** Explain, clearly and plainly. [∉, 50]

1. How do I multiply two fractions?

2. I have the numerator and the denominator of a fraction. How do I check if the fraction belongs to the open unit interval?

3. I have a positive integer. How do I check if it's prime?

4. On the plane, I have the cartesian equation of a circle, and a point. How do I decide whether or not the point lies inside the circle?

5. On the plane, I have a line and a parabola. How do I check if the line is tangent to the parabola?

6. I have two lines in three-dimensional space. How do I decide whether or not they intersect?

7. I have four points on the plane. How do I check if the points are vertices of a square?

**Exercise 2.9.** Explain concisely. [⚡, 30]

1. What is the difference between an equation and an identity?

2. What is the difference between an ordered pair and a set with two elements?

**Exercise 2.10.** I have a function between two finite sets, which are given explicitly. I can only compute the value of the function at each point of its domain, and count and compare the elements of these sets. I need explicit instructions for answering the following questions. [⚡, 50]

1. How do I check that my function is surjective?

2. How do I check that my function is injective?

**Exercise 2.11.** Answer the questions, as clearly as you can.

1. Let $A$ and $B$ be sets. Why are the sets $A$ and $(A \smallsetminus B) \cup B$ not necessarily equal? Under what condition are they equal?

2. Let $A$ and $B$ be sets. Why are the sets $(A \smallsetminus B)^2$ and $A^2 \smallsetminus B^2$ not necessarily equal? Under what condition are they equal?

3. Let $f : X \to Y$ be a function, and let $A$ be a subset of $X$. Why are the sets $A$ and $f^{-1}(f(A))$ not necessarily equal? Under what condition are they equal?

4. Let $f : X \to Y$ be a function, and let $B$ be a subset of $Y$. Why are the sets $B$ and $f(f^{-1}(B))$ not necessarily equal? Under what condition are they equal?

5. Let let $A$ and $B$ be subsets of the domain of a function $f$. Why are $f(A) \cap f(B)$ and $f(A \cap B)$ not necessarily equal? For what functions $f$ are they the same?

6. Let $\mathbf{P}(X)$ denote the power set of a set $X$. What are the elements of the set $\mathbf{P}(\mathbf{P}(\emptyset))$?

**Exercise 2.12.** Let $A, B$ be sets, and let $f, g : \mathbf{P}(A) \to \mathbf{P}(B)$. Prove that the set equation $f(x) = g(x)$ is equivalent to the equation

$$f(x) \, \Delta \, g(x) = \emptyset$$

in the sense that they have the same solution set. (Thus every set equations may be reduced to the form $F(x) = \emptyset$.)

**Exercise 2.13.** Prove that the definition

$$(a,b) \stackrel{\text{def}}{=} \{\{a\},\{a,b\}\}$$

satisfies (2.4). This shows that an ordered pair can be defined in terms of a set, so there's no need to introduce a new construct. Hence define an ordered triple in terms of sets.

**Exercise 2.14.** Consider the function that performs the prime factorization of a natural number greater than 1. What would you choose for co-domain? Explain, discussing possible representations.

**Exercise 2.15.** Find a representation for $\mathbb{Z}[x]$, the set of all polynomials with integer coefficients.

[*A polynomial is defined by the sequence of its coefficients.*]

# Chapter 3

# Logical structures

Logical structures are the mathematics of reasoning. They form the skeleton of mathematical arguments, and appear, implicitly or explicitly, in definitions, statements of facts, proofs.

Consider the following well-known result in analysis.

$$\text{Let } f : \mathbb{R} \to \mathbb{R} \text{ be a differentiable function. Then } f \text{ is continuous.} \tag{3.1}$$

This statement comprises two sentences. The first sentence does not state a fact; it's an **assumption.** We aren't even told which function $f$ is, so there is no question of this statement about $f$ being true or false. But still we use it as a basis for the rest of the argument. The second sentence does just that: it's a **deduction** of a new fact from the assumption. We seek to represent the relation between these two statements with a formal mathematical structure. The same structure should allow us to define the terms differentiability and continuity, which require the existence of certain limits at each point of the domain of $f$.

To deal with problems of this kind, we need a logical apparatus known as **predicate calculus**. This calculus has logical constants (like the numbers for ordinary calculus), logical operators (like addition or multiplication), logical functions (called predicates), and quantifiers (which act like definite integrals). The link between logical quantities and the rest of mathematics is provided by the relational operators.

## 3.1   Boolean expressions

A **boolean constant** (or boolean value) is an element of the set $\{\text{TRUE}, \text{FALSE}\}$, abbreviated to $\{\text{T}, \text{F}\}$. The boolean constants are also represented in binary notation

as 1 (TRUE) and 0 (FALSE). A **boolean expression** (or **logical expression**) is an expression that evaluates to a boolean value.

To construct logical expressions, we need **logical operators**. They are of two kinds: the **relational operators**, which convert mathematical data into logical data, and the **boolean operators**, which transform logical quantities into themselves.

## 3.1.1   Relational operators

Equalities and inequalities are the most common relational operators:

$$=  \quad  \neq  \quad  <  \quad  \leqslant  \quad  >  \quad  \geqslant . \tag{3.2}$$

The first two operators act on the elements of any set; the others act on real numbers (more generally, on the elements of an **ordered set** —see section 3.3.2).

The **relational expressions** are the simplest **mathematical sentences**, which comprise a **relational operator** and two operands. For instance, the relational expression

$$0 < 1 \tag{3.3}$$

converts a pair of real numbers into the logical constant TRUE. Any significant mathematical sentence that evaluates to TRUE is called a **theorem,** and expression (3.3) is one of the first theorems proved in analysis, which underpins all inqualities among real numbers. Expression (3.1) is also a theorem.

It is possible to state interesting facts even with simple relational expressions

$$9^3 + 10^3 = 1^3 + 12^3, \qquad \frac{355}{113} - \pi < 3 \times 10^{-7}. \tag{3.4}$$

However, it wouldn't be appropriate to call these sentences theorems, because they are very specific. More complex logical expressions such as

*The equation $10^{39} + 3 = x^2 + y^2$ has no integer solutions*

*The integer $2^{43112609} - 1$ is prime*

can still be represented as a collection of **finitely many** relational expressions. (The expression 'finitely many' denotes finiteness, without reference to actual magnitude.) The value of both expressions above turns out to be TRUE, although the verification requires a sophisticated theory and a lot of computer time[1]. In this

---

[1]The second expression represents largest known prime (as of October 2011):  see `http://primes.utm.edu/`.

course we ignore such computational difficulties (as well as deeper difficulties of theoretical nature), and we assume that any valid boolean expression is either TRUE or FALSE.

Relational expressions are an essential feature of any programming language.

```
if x<1 then
    x:=x+1
else
    x:=x-1
fi:
```

This is an example of **conditional execution**: the value of the relational expression 'x<1' determines which of two assignment statements is executed.

At the most basic level, we have the relational operators associated to sets, namely the **membership** and **subset operators** (section 2.1)

$$\in \qquad \notin \qquad \subset \qquad \not\subset . \tag{3.5}$$

For example, in section 6.1 we prove that the relational expression

$$\sqrt{2} \notin \mathbb{Q} \tag{3.6}$$

is TRUE, namely, that the square root of 2 is not a rational number.

There are countless relational operators in mathematics: the **divisibility operator** '|' and the **congruence operator** '≡' in arithmetic (see section 2.1.2), the **isomorphism operator** '≅' in algebra, the **orthogonality operator** '⊥' in geometry, etc.

The symbol '≈', used in mathematical physics, does not represent a relational operator, because expressions such as $\pi \approx 3.14$ cannot be assigned unequivocally a value TRUE or FALSE.

## 3.1.2 Boolean operators

Boolean expressions may be constructed from boolean constants and relational expressions by means of **boolean operators**. This process is analogous to the construction of arithmetical expressions from arithmetical constants (numbers) and operators ($+$, $-$, etc.).

The basic boolean operators are NOT, and AND, represented symbolically as $\neg$ and $\wedge$, respectively. All other operators can be expressed in terms of them. The operator NOT is **unary,** that is, it takes just one boolean operand and produces a

boolean result. In this respect the operator NOT resembles the unary arithmetical operator '$-$', which changes the sign of a number. The operator AND is **binary** —it acts on two operands  and produces a boolean result. Think of it as a kind of 'multiplication'.

The following **truth table** defines the operators NOT and AND, by specifying its action on all possible choices of operands

$$
\begin{array}{c|c}
P & \neg P \\
\hline
T & F \\
F & T
\end{array}
\qquad
\begin{array}{c|c|c}
P & Q & P \wedge Q \\
\hline
T & T & T \\
T & F & F \\
F & T & F \\
F & F & F
\end{array}
\tag{3.7}
$$

The expression $\neg P$ is called the **negation** of $P$.  Negating relational expressions is straightforward, since we already have all relevant symbols —see (3.2) and (3.5). Thus

$$
\neg(x < y) = (x \geqslant y), \qquad \neg(x \in A) = (x \notin A), \qquad \neg(A \not\subset B) = A \subset B.
$$

If $P$ and $Q$ are boolean expression, then $P \wedge Q$ is called a **compound expression,** or compound statement. In a compound statement, the operator AND  may appear implicitly. For instance, the expression $0 < x < 1$ is a compound boolean expressions, which is written in full as $(x > 0) \wedge (x < 1)$.

Other binary operators may be constructed from the above two. The most commonly used are OR (represented by the symbol $\vee$) and the **implication operator** $\Rightarrow$. We define these operators with truth tables, although they could also be defined in terms of NOT and AND (see remarks following theorem 3.1.2, below).

$$
\begin{array}{c|c|c}
P & Q & P \vee Q \\
\hline
T & T & T \\
T & F & T \\
F & T & T \\
F & F & F
\end{array}
\qquad
\begin{array}{c|c|c}
P & Q & P \Rightarrow Q \\
\hline
T & T & T \\
T & F & F \\
F & T & T \\
F & F & T
\end{array}
\tag{3.8}
$$

We see from the table that the operator OR is **inclusive,** namely it always includes the possibility that both operands are true, which is not necessarily the meaning attributed to the conjunction 'or' in common English usage.

The expression

$$
P \Rightarrow Q
\tag{3.9}
$$

where $P, Q$ are boolean expressions is called an **implication,** which reads

| | | |
|---|---|---|
| *P implies Q* | *Q follows from P* | *if P, then Q* |
| *P only if Q* | *P is sufficient for Q.* | |

The statement $P$ is the **hypothesis,** while $Q$ is the **conclusion.** Because many mathematical statements have this form (the statement (3.1), for example), we shall study this operator thoroughly.

As with arithmetical operators, the order of evaluation of boolean operators may be altered by means of parentheses

$$(\text{F} \wedge \text{T}) \Rightarrow \text{T} \qquad \text{F} \wedge (\text{T} \Rightarrow \text{T}).$$

The expressions above evaluate to TRUE and FALSE, respectively. In absence of parentheses, the implication operator binds more tightly than $\wedge$ and $\vee$, while the latter have the same precedence. If there is more than one operator, the evaluation proceeds from left to right. Thus the expression $P \vee Q \Rightarrow R \wedge S$ is evaluated as $((P \vee Q) \Rightarrow R) \wedge S$. However, these rules are not as established as those for arithmetical operators, and we parentheses should be used whenever there is a potential ambiguity. A relational operator binds more tightly than any boolean operator, thereby removing any ambiguity. Thus the expression $1 > 2 \Rightarrow 2 > 1$ can only be interpreted as $(1 > 2) \Rightarrow (2 > 1)$, because $1 > (2 \Rightarrow 2) > 1$ is syntactically incorrect (the operands of $\Rightarrow$ must be boolean). Redundant parentheses may still be used to improve readability.

From the truth table (3.8) we see that the expression FALSE $\Rightarrow Q$ is true for any value of $Q$. This definition, which is perhaps unexpected, is not only essential to the theory, but also in agreement with the common usage of implications. Consider the sentences

*If 4 is a prime number, then I am a martian.*
*If I win the lottery, then I'll buy you a Ferrari.*

The first sentence, which makes sense in English, is also true mathematically. As to the second, there are four possible scenarios, and the statement is false in only one case (I win the lottery but I don't buy you a Ferrari), in agreement with our convention.

We remark that $P \Rightarrow Q$ is different from $Q \Rightarrow P$, that is, the operator $\Rightarrow$ (unlike $\wedge$ and $\vee$) is **non-commutative.** For this reason, one also define the operator $\Leftarrow$, given in terms of $\Rightarrow$ by interchanging operands, namely

$$P \Leftarrow Q \stackrel{\text{def}}{=} Q \Rightarrow P.$$

The expression on the left is called the **converse** of the implication (3.9), and it reads

      *P is implied by Q*           *P follows from Q*

      *P if Q*                  *P is necessary for Q.*

The boolean value of an implication is unrelated to the value of the converse implication. In the following example, the former is false and the latter is true.

$$(x^2 = 25) \Rightarrow (x = -5) \qquad (x^2 = 25) \Leftarrow (x = -5).$$

Inappropriate reversal of an implication is a common mistake in proofs —see section 9.3.

The double-headed arrow

$$P \Leftrightarrow Q \tag{3.10}$$

is the **equivalence operator**, which is defined in terms of the operators $\Rightarrow$ and $\wedge$ as follows

$$P \Leftrightarrow Q \stackrel{\text{def}}{=} (P \Rightarrow Q) \wedge (Q \Rightarrow P). \tag{3.11}$$

One verifies that its truth table is

| $P$ | $Q$ | $P \Leftrightarrow Q$ |
|---|---|---|
| T | T | T |
| T | F | F |
| F | T | F |
| F | F | T |

The meaning of $\Leftrightarrow$ results from the conjunction of $\Rightarrow$ and $\Leftarrow$; so the expression (3.10) is read out loud as

      *P implies and is implied by Q*        *P is equivalent to Q*

      *P if and only if Q*              *P is necessary and sufficient for Q*

The awkward expression *'if and only if'* (abbreviated 'iff') is much in use: thus the statement

$$A \subset B \Leftrightarrow A \smallsetminus B = \emptyset \qquad A \subset B \text{ iff } A \smallsetminus B = \emptyset.$$

may be read out loud as

*A is a subset of B if and only if the set difference of A and B is empty.*

One description of $P \Leftrightarrow Q$ is that $P$ and $Q$ are **equivalent statements** —they are both true or both false. In Boolean expressions, equivalence replaces equality. So, of the two expressions

$$\neg(\neg P) \Leftrightarrow P \qquad \neg(\neg P) = P$$

only the left one is —strictly speaking— syntactically correct. On the other hand, the symbol '=' is so embedded in mathematical writing, that it hardly seems sensible to ban it from predicate calculus. In addition, excessive usage of $\Leftrightarrow$ may make an expression difficult to read. For these reasons, the equal sign may also be used to mean equivalence, as long as this does not lead to confusion.

The **contrapositive** of the implication (3.9) is the implication

$$\neg P \Leftarrow \neg Q. \tag{3.12}$$

In other words, the contrapositive of an implication is constructed by reversing the implication *and* negating the operands. Great care must be exercised in distinguishing between the direct implication, its converse, and its contrapositive.

| | | |
|---|---|---|
| DIRECT : | *If x is a multiple of 4, then x is even.* | (TRUE) |
| CONVERSE : | *If x is even, then x is a multiple of 4.* | (FALSE) |
| CONTRAPOSITIVE : | *If x is odd, then x is not a multiple of 4.* | (TRUE) |

While the value of an implication and of its converse are unrelated, *every implication is equivalent to its contrapositive* —they are both true or false for any choice of the operands. This is the identity (iv) of the following theorem.

**Theorem.** *For all $P, Q \in \{\text{TRUE}, \text{FALSE}\}$, the following holds*

$$\begin{aligned}
(i) && \neg(P \vee Q) &\Leftrightarrow (\neg P \wedge \neg Q) \\
(ii) && \neg(P \wedge Q) &\Leftrightarrow (\neg P \vee \neg Q) \\
(iii) && P \Rightarrow Q &\Leftrightarrow (\neg P \vee Q) \\
(iv) && P \Rightarrow Q &\Leftrightarrow (\neg P \Leftarrow \neg Q).
\end{aligned}$$

PROOF. The proof consists of evaluating each side of these equalities for all possible values of $P, Q$. We prove (*iii*). The other proofs are left as an exercise. The left-hand

side of (*iii*) was given in (3.8). We compute the right-hand side explicitly.

| $P$ | $Q$ | $\neg P$ | $\neg P \vee Q$ | $P \Rightarrow Q$ |
|-----|-----|----------|-----------------|-------------------|
| T   | T   | F        | T               | T                 |
| T   | F   | F        | F               | F                 |
| F   | T   | T        | T               | T                 |
| F   | F   | T        | T               | T                 |

We see that the left-hand side and the right-hand side of (*iii*) are equal.    □

The statements (*i*) and (*ii*) are known as *De Morgan's laws*[2].    Using theorem
3.1.2, one can express the operators $\vee, \Rightarrow, \Leftrightarrow$ in terms of $\vee$ and $\wedge$ (see exercises).

EXAMPLE.  Comment on the following expression  **[∉]**

$$P \vee Q := \neg((\neg Q) \wedge (\neg P)).$$

*This is a boolean expression, which defines the operator* OR *in terms
of the operators* NOT *and* AND *. This definition shows that the operator*
OR *need not be defined with a truth table.*


## 3.2   Predicates

A **predicate** (or **boolean function**) $\mathscr{P}$ is a function assuming boolean values. This
simple idea has far-reaching applications. If the domain of $\mathscr{P}$ is $X$, then we speak
of a **predicate on** (or **over**) $X$.  The following function

$$\mathscr{P} : \mathbb{Z} \to \{\text{T},\text{F}\} \qquad x \mapsto 7 \mid x \tag{3.13}$$

is a predicate on the integers; we verify that $\mathscr{P}(22) = \text{F}$, $\mathscr{P}(-91) = \text{T}$.

Let $\mathscr{P} : X \to \{\text{T},\text{F}\}$ be a predicate on $X$. There is a distinguished subset $A$ of $X$,
determined by $\mathscr{P}$ via the following Zermelo definition

$$A := \{x \in X \: : \: \mathscr{P}(x)\}. \tag{3.14}$$

With reference to equation (2.7) and the discussion that follows, we see that the
expression $\mathscr{P}(x)$ which means '*x* has property *P*', is the definiens of the set $A$. Ne
now have the opportunity to express the Zermelo definition of a set in the language
of functions

$$\{x \in X \: : \: \mathscr{P}(x)\} \stackrel{\text{def}}{=} \mathscr{P}^{-1}(\{\text{T}\}).$$

---

[2]Augustus De Morgan, British mathematician and logician (1806–1871).

We see that the set $A$ in (3.14) is just an inverse image of a set under function (cf. definition (2.15), page 25.)

Conversely, let $X$ be a set and let $A$ be a subset of $X$. The predicate

$$\mathscr{P}_A : X \to \{\text{T},\text{F}\} \qquad x \mapsto x \in A \qquad (3.15)$$

is called the **characteristic function** of $A$ (in $X$). (Explicitly, $\mathscr{P}(x)$ is equal to T if $x$ belongs to $A$, and to F otherwise.) For example, the function (3.13), is the characteristic function of $7\mathbb{Z}$, the set of integer multiples of 7.

So to every predicate on a set $X$ we associate a unique subset of $X$, and vice-versa. We say that there is a **bi-unique correspondence** between these two classes of objects, established by expressions (3.15) and (3.14).

Boolean expressions are manipulated using boolean operators; sets are manipulated using set operations. The following theorem establishes correspondences between the two classes of objects. (In section 5.7 we'll write a short essay about this theorem.)

**Theorem.** *Let $X$ be a set, let $A, B \subseteq X$, and let $\mathscr{P}_A$, $\mathscr{P}_B$ be the corresponding characteristic functions. The following holds (the prime denotes taking the complement)*

$$
\begin{array}{rrcl}
(i) & \neg\, \mathscr{P}_A & = & \mathscr{P}_{A'} \\
(ii) & \mathscr{P}_A \wedge \mathscr{P}_B & = & \mathscr{P}_{A\cap B} \\
(iii) & \mathscr{P}_A \vee \mathscr{P}_B & = & \mathscr{P}_{A\cup B} \\
(iv) & \mathscr{P}_A \Rightarrow \mathscr{P}_B & = & \mathscr{P}_{(A\smallsetminus B)'} \\
(v) & \mathscr{P}_A \Leftrightarrow \mathscr{P}_B & = & \mathscr{P}_{(A\cap B)\cup(A\cup B)'}.
\end{array}
$$

PROOF. To prove (*i*) we note that the function $x \mapsto \neg\mathscr{P}_A(x)$ evaluates to TRUE if $x \notin A$ and to FALSE otherwise. However, from the definition of complement of a set, we have $x \notin A \Leftrightarrow x \in A'$.

Next we prove (*iv*). We'll prove instead

$$\neg(\mathscr{P}_A \Rightarrow \mathscr{P}_B) = \mathscr{P}_{(A\smallsetminus B)}$$

which, together with (i), gives us (iv). Let $\mathscr{P}_A := (x \in A)$ and let $\mathscr{P}_B := (x \in B)$. Then, from the truth table (3.8) for the operator $\Rightarrow$, we find that $\neg(\mathscr{P}_A \Rightarrow \mathscr{P}_B)$ is TRUE precisely when $\mathscr{P}_A(x)$ is TRUE and $\mathscr{P}_B(x)$ is FALSE. This means that

$$(x \in A) \wedge (x \notin B)$$

but this is just the definition of the characteristic function of the set $A \smallsetminus B$, as desired.

The proof of (*ii*), (*iii*), (*v*) is left as an exercise.    □

We illustrate the significance of this theorem with some examples.

EXAMPLE. Let $a$ and $b$ be real numbers. The predicates $x \mapsto (x \geqslant a)$ and $x \mapsto (x < b)$ (over $\mathbb{R}$) are the characteristic functions of two rays, one of which without endpoint. According to theorem 3.2, part (*ii*), the predicate $x \mapsto ((x \geqslant a) \wedge (x < b))$ is the characteristic function of the intersection of these rays. Depending on whether $a < b$, or $a \geqslant b$, this intersection is the half-open interval $[a, b)$ or the empty set.

EXAMPLE. If the sets $A$ and $B$ are disjoint, then $A \cap B$ is empty, and from part (iv) of the theorem, we obtain

$$(\mathscr{P}_A \Rightarrow \mathscr{P}_B) = \mathscr{P}_{\emptyset'} = \mathscr{P}_X.$$

For example, we have $4\mathbb{Z} \subset 2\mathbb{Z}$, and hence, over the integers

$$(\mathscr{P}_{4\mathbb{Z}} \Rightarrow \mathscr{P}_{2\mathbb{Z}}) = \mathscr{P}_{\mathbb{Z}}.$$

This functional identity encodes the truth of the statement that every multiple of four is even.

EXAMPLE. Let $X, Y$ sets, and let $f, g : X \to Y$ be functions. An **equation** (on $X$) is a predicate of the type

$$\mathscr{P} : X \to \{\mathrm{T}, \mathrm{F}\} \qquad x \mapsto (f(x) = g(x)).$$

In the language of predicates, an equation is the characteristic function of its own solution set $\{x \in X : \mathscr{P}(x)\}$ (cf. section 2.5). The system of two equations

$$f_1(x) = g_1(x) \qquad f_2(x) = g_2(x)$$

defined over the same set corresponds to the predicate $x \mapsto \mathscr{P}_1(x) \wedge \mathscr{P}_2(x)$. From theorem 3.2, part (*ii*), it follows that the solution set of a system of two equations is the intersection of the solution sets of the individual equations.

## 3.3   Quantifiers

Theorems (3.3) and (3.6) of section 3.1.1, are very specific. By contrast, the statement (3.1) is general, in that it refers to a family of objects, rather than an individual

object. The statement of general facts requires two special symbols, $\forall$ and $\exists$, called **quantifiers**.

| | | | | |
|---|---|---|---|---|
| **Universal quantifier:** | $\forall$ | *'for all'* | *'for any choice of'* | *'given any'* |
| **Existential quantifier:** | $\exists$ | *'there exists'* | *'for some'* | *'we can find'* |

Formulating interesting mathematical statements would be impossible without these innocuous-looking symbols. Indeed the number of quantifiers required to define a concept may be taken as an indicator of the concept's logical depth. For instance, the definition of continuity in analysis requires at three quantifiers —see section 4.5.

A quantified expression has the following syntax

$$\forall x \in X, \mathscr{P}(x) \qquad \text{For all } x \text{ in } X, \mathscr{P}(x) \tag{3.16}$$

$$\exists x \in X, \mathscr{P}(x) \qquad \text{There exists } x \text{ in } X, \text{ such that } \mathscr{P}(x). \tag{3.17}$$

where $X$ is a set, and $\mathscr{P}$ is a predicate over $X$. The quantifier is followed by the symbol being quantified, and then by the membership of the latter to a set. The predicate is required; the expressions $\forall x \in X$ and $\exists x \in X$ are incomplete and have no meaning.

The meaning (boolean value) of the expressions (3.16) and (3.17) is given by the following equivalences

$$\forall x \in X, \mathscr{P}(x) \overset{\text{def}}{\Leftrightarrow} \mathscr{P}(X) = \{\text{TRUE}\}$$
$$\exists x \in X, \mathscr{P}(x) \overset{\text{def}}{\Leftrightarrow} \mathscr{P}(X) \neq \{\text{FALSE}\} \tag{3.18}$$

which *define* the quantified expressions on the left in terms of the unquantified expressions on the right.

As an example, consider the following sentences

$$\forall x \in \mathbb{Z}, \; x < 0 \qquad \text{All integers are negative}$$
$$\exists x \in \mathbb{Z}, \; x < 0 \qquad \text{There is a negative integer}$$

which are FALSE and TRUE, respectively. The English expressions are a synthesis of a literal translation of symbols into words

*Given any integer $x$, $x$ is negative*
*We can find an integer $x$ such that $x$ is negative.*

Significantly, in both cases the quantified variable $x$ has become silent —no explicit reference to it remains in the sentence. This is a general phenomenon, as we shall see below.

The acronym 's.t.' (for 'such that') is sometimes inserted in expressions involving $\exists$

$$\exists x \in X \quad \text{s.t.} \quad \mathscr{P}(x).$$

This abbreviation may improve the readability of a mathematical sentence, even if it plays no formal role. If the ambient set is clear from the context, reference to it may be omitted, particularly in conjunction with an inequality that restricts the variable's range. For example

$$\forall n > 3, \ n! > 2^n \qquad \text{means} \qquad \forall n \in \mathbb{N} \smallsetminus \{1,2,3\}, \ n! > 2^n.$$

EXAMPLE. Using a quantifier, we can *define* the subset operator in terms of the membership operator

$$A \subset B \overset{\text{def}}{\Leftrightarrow} \forall x \in A x \in B. \tag{3.19}$$

Quantifiers are powerful symbols, particularly when combined together. To give a taste of the limitless possibilities, let us consider the predicate

$$\mathscr{L}(x,y) := (\text{'}x \text{ loves } y\text{'}).$$

This is a function of two variables, which we define over the cartesian product $X = G \times G$, where $G$ is a set of people. So the lovers and the loved ones are extracted from the same set.

Now choose $g \in G$, say, $g =$ George. Then 'everybody loves George' is a boolean expression, which is described mathematically as

> *For all elements $x$ of the set $G$, the value of the expression $\mathscr{L}(x,g)$ is* TRUE.

Using a quantifier, this expression translates into symbols as follows

$$\forall x \in G, \ \mathscr{L}(x,g).$$

This expression may be true or false (depending on how charming George is). The amount of computation required to evaluate the above expression (the number of evaluations of the function $\mathscr{L}$ could be as high as $\#G$) gives an idea of the power of

a quantified expression. The following examples illustrate some possibilities, with one and two quantifiers.

| | words | symbols |
|---|---|---|
| 1 | *everybody loves George* | $\forall x \in G, \; \mathcal{L}(x, g)$ |
| 2 | *somebody loves George* | $\exists x \in G \smallsetminus \{g\}, \; \mathcal{L}(x, g)$ |
| 3 | *everybody loves himself* | $\forall x \in G, \; \mathcal{L}(x, x)$ |
| 4 | *George loves nobody* | $\forall x \in G, \; \neg \mathcal{L}(g, x)$ |
| 5 | *George is in love* | $\exists x \in G \smallsetminus \{g\}, \; \mathcal{L}(g, x)$ |
| 6 | *everybody loves somebody* | $\forall x \in G, \; \exists y \in G, \; \mathcal{L}(x, y)$ |
| 7 | *somebody loves everybody* | $\exists x \in G, \; \forall y \in G, \; \mathcal{L}(x, y)$ |
| 8 | *somebody is loved by everybody* | $\exists y \in G, \; \forall x \in G, \; \mathcal{L}(x, y)$ |

In examples 2 and 5, the indeterminate $x$ belongs to the set $G \smallsetminus \{g\}$, to ensure that $x$ is distinct from $g$. It is possible to transfer this constraint to the predicate. Thus example 2 is equivalent to

$$\exists x \in G, \; (x \neq g) \wedge \mathcal{L}(x, g).$$

Examples 6 and 8 illustrate an important fact: exchanging the order of quantifiers may alter completely the meaning of an expression. These expressions, which involve two quantifiers, have implied parentheses, which establish the order of evaluation of the various sub-expressions. For instance expression 6 is written in full as

$$\forall x \in G, \; (\exists y \in G, \; \mathcal{L}(x, y))$$

For this expression to be valid, the quantity in parentheses must define a predicate over $G$. We now show that this is indeed the case. Let $x \in G$ be given. Then $\{x\} \times G$ is a subset of $G \times G$, and therefore we can form its image under $\mathcal{L}$, obtaining a subset of $\{\text{T}, \text{F}\}$. Our predicate is now defined using (3.18)

$$\exists y \in G, \; \mathcal{L}(x, y) \Leftrightarrow \mathcal{L}(\{x\} \times G) \neq \{\text{FALSE}\}.$$

(Think about it.)

Thus quantifiers are **operators,** which act on boolean functions producing new functions. Their effect is to reduce the number of variables by one, a process that calls to mind definite integration. One should think of the two expressions

$$\forall x \in X \qquad \int_X dx$$

as being structurally similar.  The operator on the left acts on predicates over the set $X$, namely on functions $\mathscr{P} : X \to \{T, F\}$.  The operator on the right acts, say, on functions $F : X \to \mathbb{R}$, which are integrable over a subset $X$ of the real line. (To fix ideas, think of $X$ as an interval.)  Inserting the appropriate functions in each expression

$$\forall x \in X, \ \ \mathscr{P}(x) \qquad\qquad \int_X dx\, F(x)$$

we obtain a boolean constant (TRUE or FALSE) on the left, and a numerical constant (a number) on the right.  Now suppose that both $\mathscr{P}$ and $F$ are functions of two variables $x$ and $y$, defined over the cartesian product $X \times Y$ of two sets.  Then both expressions

$$\forall x \in X, \ \ \mathscr{P}(x, y) \qquad\qquad \int_X dx\, F(x, y)$$

produce functions of $y$, which is the variable that is not quantified in one case, and not integrated over in the other.  If we quantify/integrate with respect to both variables,

$$\forall x \in X, \ \ \exists y \in Y, \ \ \mathscr{P}(x, y) \qquad\qquad \int_X dx \int_Y dy\, F(x, y)$$

we obtain again constants.

Returning to the examples above, we note that $\mathscr{L}(x, g)$ is a predicate in one variable ($g$ is fixed), and so is $\mathscr{L}(x, x)$.  Hence $\exists x, \mathscr{L}(x, g)$, and $\forall x, \mathscr{L}(x, x)$ are constants. On the other hand, any expressions in which the numbers of quantifiers is smaller than the number of arguments in the predicate is a predicate with fewer arguments, as the following examples illustrate.

| words | symbols |
|---|---|
| *x is in love* | $\exists y \in G \smallsetminus \{x\}, \ \mathscr{L}(x, y)$ |
| *x is loved* | $\exists y \in G \smallsetminus \{x\}, \ \mathscr{L}(y, x)$ |
| *x is a hippy* | $\forall y \in G, \ \mathscr{L}(x, y)$ |
| *x is selfish* | $\mathscr{L}(x, x) \wedge [\forall y \in G \smallsetminus \{x\}, \ \neg\mathscr{L}(x, y)]$ |
| *x is a lover of George* | $x \neq g \wedge \mathscr{L}(x, g) \wedge \mathscr{L}(g, x)$ |

Now, each predicate is the characteristic function of a subset of $G$, which we construct with a Zermelo definition.

| words | symbols |
|-------|---------|
| *the people in love* | $\{x \in G : \exists y \in G \smallsetminus \{x\},\ \mathscr{L}(x,y)\}$ |
| *the loved ones* | $\{y \in G : \exists x \in G \smallsetminus \{y\},\ \mathscr{L}(x,y)\}$ |
| *the hippies* | $\{x \in G : \forall y \in G,\ \mathscr{L}(x,y)\}$ |
| *the selfish people* | $\{x \in G : \mathscr{L}(x,x) \wedge \forall y \in G \smallsetminus \{x\},\ \neg\mathscr{L}(x,y)\}$ |
| *George's lovers* | $\{x \in G : x \neq g \wedge \mathscr{L}(x,g) \wedge \mathscr{L}(g,x)\}.$ |

Let us consider a quantified expression of the form

$$\forall x \in X,\ \exists y(x) \in Y,\ \mathscr{P}(x,y) \tag{3.20}$$

where $X$ and $Y$ are sets and $\mathscr{P}$ is a predicate over $X \times Y$. The notation makes it clear that the choice of $y$ depends on $x$, in general. The sentence "*everybody loves somebody*" considered above is of this form, with $X = Y = G$, a set of people, and $\mathscr{P}(x,y) = $ '$x$ loves $y$'. The need for an explicit dependence of $y$ on $x$ is obvious in this case. Let us now change the meaning of the symbols, while keeping the same structure. Letting $X = Y = \mathbb{N}$, and $\mathscr{P}(x,y) = (x < y)$, we obtain the statement

*Given any integer, one can find a larger integer*

or, better,

*There is no greatest integer.*

This statement is clearly unrelated to the previous one, yet their formal structure is the same: a universal and an existential quantifier acting on a predicate in two variables.

It's useful to think of the interplay between universal and existential quantifiers as a representation of an adversarial system —like a court of law— based on the following rules of engagement:

| | | |
|---|---|---|
| $\forall$ | *given any* | (my opponent's move) |
| $\exists$ | *I can find* | (my move). |

The quantifiers create a tension between the two contenders, and I should expect my opponent to challenge me. Accordingly, we rewrite our statement more emphatically

*Given any integer, no matter how large, one can always find a larger integer.*

The expressions 'no matter how large', and 'always' are inessential to the claim, but they illustrate the dynamics of the process. This aspect of predicate calculus will be considered again in chapter 6, when we deal with proofs. We shall see that the formal structure of a statement shapes the structure of a proof.

EXAMPLE. The **Archimedean property** of the real numbers states that

> *The integer multiples of any positive quantity can be made arbitrarily large.*

The symbolic version of the statement requires three quantifiers.

$$\forall x \in \mathbb{R}^+, \ \forall y \in \mathbb{R}, \ \exists n \in \mathbb{N}, \ nx > y.$$

In this expression $x$ is the positive quantity in question, $y$ is the (large) quantity we want to exceed, and $n$ is the multiple of $x$ needed to achieve this. Note that $n = n(x,y)$. It seems easier to define the Archimedean property with words than with symbols.

## 3.3.1    Quantifiers and functions

A statement about a function will invariably refer to the elements of the function's domain or co-domain. Quantifiers may also act on these elements. We begin with injectivity and surjectivity (see section 2.2). A function is **injective** if it maps distinct points to distinct points. This concise statement is spelled out as follows

> *Given any two points in the domain of the function, if they are distinct, then so are their images under the function.*

The presence of 'given any' and 'if ... then' tells us that there is a universal quantifier and an implication operator. Let $f : A \rightarrow B$. A literal translation of the above expression into symbols is

$$\forall x, y \in A, \ (x \neq y) \Rightarrow (f(x) \neq f(y))$$

where $\forall x, y \in A$ is a shorthand for $\forall x \in A, \ \forall y \in A$. Replacing the implication by its contrapositive (theorem 3.1.2 (iv), page 55) we obtain a neater expression

$$\forall x, y \in A, \ (f(x) = f(y)) \Rightarrow (x = y).$$

The original definition ('distinct points have distinct images') restricts the ambient set to pairs of distinct point. This can be done symbolically:

$$\forall x \in A, \ \forall y \in A \smallsetminus \{x\}, \ f(x) \neq f(y).$$

The predicate is now simpler, but the description of the underlying sets has become more complicated. It is always possible to transform an implication in this way, by absorbing the hypothesis into the ambient set, see exercise 3.

A function $f$ is **surjective** if every point in the co-domain is the image of some point in the domain, that is, if

> *Given any point in the co-domain of the function, we can find a point in the domain which maps to it.*

This time we need both quantifiers

$$\forall y \in B, \ \exists x \in A, \ f(x) = y.$$

This expression is of the form (3.20).

In section 2.4) we noted that a sequence of elements of a set $A$ may be interpreted as a function

$$a : \mathbb{N} \to A \qquad k \mapsto a_k.$$

Characterising sequences is then analogous to characterising functions.

For example, consider the set $\mathbb{Z}^{\mathbb{N}}$ of all integer sequences (this notation was developed in section 2.4). To isolate sequences with certain properties —a subset of $\mathbb{Z}^{\mathbb{N}}$— we use a Zermelo definition. For example, the set

$$\{a \in \mathbb{Z}^{\mathbb{N}} : a_1 \in 2\mathbb{Z}\}$$

consists of all integer sequences whose first term is even. If we apply the quantifier $\forall$ to the subscript (which is the variable in our functions), we'll be able to deal with all terms of the sequences at once. So the set of integer sequences with only even terms is given by

$$\{a \in \mathbb{Z}^{\mathbb{N}} : \forall k \in \mathbb{N}, \ a_k \in 2\mathbb{Z}\} = (2\mathbb{Z})^{\mathbb{N}}. \tag{3.21}$$

There seems to be something wrong here. The quantifier 'integrates out' the indeterminate $k$: does this mean that the predicate is just a constant? No, because here the relevant variable is actually $a = (a_k)$, an element of the ambient set $\mathbb{Z}^{\mathbb{N}}$. Indeed, if we write (3.21) as $\{a \in \mathbb{Z}^{\mathbb{N}} : \mathscr{P}(a)\}$, then the predicate is given explicitly by

$$\mathscr{P} : \mathbb{Z}^{\mathbb{N}} \to \{\text{T}, \text{F}\} \qquad a \mapsto (\forall k \in \mathbb{N}, \ a_k \in 2\mathbb{Z}).$$

Likewise, the set

$$\{a \in \mathbb{Z}^{\mathbb{N}} : \exists k \in \mathbb{N}, \ a_k \in 2\mathbb{Z}\}$$

is

> *The set of integer sequences with one even term.*

The above expression reflects the very literal interpretation of what's written, which is what mathematicians (and lawyers) are known for. A mathematical geek would be entertained by the fact that the statement '*In London there is one tube station*', is true; a normal person would instead perceive this as a puzzling understatement ('Where do they go from there?'). So, for an effective delivery, a characterisation of the type

> *The set of integer sequences with at least one even term.*

is preferable to the one given before. The qualifier 'at least', while, strictly speaking, redundant, helps the reader note an essential point.

The set

$$\{a \in \mathbb{Z}^{\mathbb{N}} : \forall n \in \mathbb{N}, \ n > 2 \Rightarrow 2|a_n\}$$

is described as

> *The set of integer sequences whose terms, after the second one, are even.*

Even if information is given about the first two terms, this statement could be misinterpreted as meaning that these terms are not even. A more helpful description of this set is

> *The set of integer sequences whose terms are all even, with the possible exception of the first two terms.*

From these examples, we begin to learn that much can be done to help the reader understand a mathematical text. The question of effectiveness in mathematical writing will be considered systematically in chapter 5.

## 3.3.2   Quantifiers and relations

Let $X$ and $Y$ be sets. A **relation** $\mathscr{R}$ on $X \times Y$ is a predicate over $X \times Y$. If $X = Y$, we speak of a **relation on** $X$. If $\mathscr{R}$ is a relation, it is customary to write $x\mathscr{R}y$ to mean $\mathscr{R}(x, y)$. The expression $x\mathscr{R}y$ is called a **relational expression**. All relational

expressions introduced a the beginning of this chapter are of this form. Thus the membership operator $\in$ defines a relation on $X \times \mathbf{P}(X)$, where $X$ is some ambient set.

(A relation on a set $X$ is sometimes defined as a **subset** of $X^2$, rather then a predicate over $X^2$. In this sense, the set $\{(1,1),(2,1)\}$ is a relation on $\{1,2\}$. The correspondence between sets and predicates described in section 3.2 clarifies the connection between the two constructs.

A relation $\mathcal{R}$ on a set $X$ is called an **equivalence relation** if it satisfies the following properties

$$\forall x \in X, \quad x\mathcal{R}x \qquad\qquad \textbf{reflexivity}$$
$$\forall x,y \in X, \quad x\mathcal{R}y \Rightarrow y\mathcal{R}x \qquad\qquad \textbf{symmetry}$$
$$\forall x,y,z \in X, \quad (x\mathcal{R}y \wedge y\mathcal{R}z) \Rightarrow x\mathcal{R}z \qquad \textbf{transitivity.}$$

EXAMPLE. The relational operator '=' defines an equivalence relation on any set. This is the **trivial** equivalence.

EXAMPLE. For any natural number $m$, the congruence relation $x \equiv y \,(\mathrm{mod}\ m)$ defined in section 2.1.2 is an equivalence relation on $\mathbb{Z}$.

EXAMPLE. The relation '$\sim$' on $\mathbb{N} \times \mathbb{N}$, defined by $(m,n) \sim (j,k)$ if $m + k = n + j$, is an equivalence relation. By interpreting the pair $(m,n)$ as the quantity $z = m - n$, we see that equivalent pairs correspond to the same value of $z$. With this device, one can construct integers from pairs of natural numbers. More precisely, every integer is an infinite set of equivalent natural numbers. This is the **abstract definition** of the integers. The virtue of this construction is that it requires only addition in $\mathbb{N}$, not subtraction.

EXAMPLE. The relation '$\sim$' on $\mathbb{Z} \times (\mathbb{Z} \smallsetminus \{0\})$, defined by $(m,n) \sim (j,k)$ if $mk = nj$, is an equivalence relation. By interpreting the pair $(m,n)$ as $r = m/n$, we see that equivalent pairs correspond to the same value of $r$. With this device, one can define a rational number as an infinite collection of equivalent pairs of integers. This is the **abstract definition** of the rational numbers, which requires only multiplication in $\mathbb{Z}$, not division.

A relation $\mathcal{R}$ on a set $X$ is called a **partial ordering** if it satisfies the following properties

$$\forall x \in X, \quad x\mathcal{R}x \qquad\qquad \textbf{reflexivity}$$
$$\forall x,y \in X, \quad (x\mathcal{R}y \wedge y\mathcal{R}x) \Rightarrow x = y \qquad \textbf{anti-symmetry}$$
$$\forall x,y,z \in X, \quad (x\mathcal{R}y \wedge y\mathcal{R}z) \Rightarrow x\mathcal{R}z \qquad \textbf{transitivity.}$$

A set is **partially ordered** if a partial ordering is defined on it. A partial ordering is usually denoted by the symbol '$\leqslant$'. So we write $x \leqslant y$ instead of $x \mathscr{R} y$.

EXAMPLE. The relational operator $\leqslant$ defines a partial ordering in $\mathbb{N}$, $\mathbb{Z}$, $\mathbb{Q}$, and $\mathbb{R}$, but not in $\mathbb{C}$.

EXAMPLE. The set $\mathbf{P}(X)$ of all subsets of a set $X$ is partially ordered by set inclusion, whereby $\leqslant$ means $\subset$.

A partially ordered set $X$ is said to be **ordered**, if all pairs of elements of $X$ are **comparable**, meaning that we either have $x \leqslant y$ or $y \leqslant x$. The real line is an ordered set; the power set $\mathbf{P}(X)$ of a set $X$, which is partially ordered by set inclusion, is not ordered.

An ordered set $X$ is said to be **well-ordered** if any non-empty subset $A \subset X$ has a smallest element. The symbolic definition requires three quantifiers

$$\forall A \in \mathbf{P}(X) \smallsetminus \emptyset, \ \exists a \in A, \ \forall x \in A, \ a \leqslant x.$$

EXAMPLE. Any finite ordered set is well-ordered. The closed unit interval $[0,1]$ is ordered but not well-ordered, because the subset $(0,1]$ has no smallest element. The natural numbers are well-ordered. This property forms the basis of the principle of induction, which we consider in chapter 7.

## 3.4   Existence statements

An **existence statement** is a boolean expression with a leading existential quantifier $\exists$. One must keep in mind that the quantifier may be hidden; thus the sentence

*The number* $\cos(\pi/3)$ *is rational*

is an existence statement. The hidden quantifier becomes visible if we write

*For some rational number r, we have* $\cos(\pi/3) = r$

or, in symbols,

$$\exists r \in \mathbb{Q}, \ r = \cos(\pi/3).$$

This statement is weaker than the identity $\cos(\pi/3) = 1/2$, because it does not require us to specify which rational number our expression is equal to.

An existential quantifier is hidden in the definition of divisibility of integers. We say that $a$ divides $b$ if there exists an integer $m$ such that $am = b$. In symbols

$$a|b \ \Leftrightarrow \ \exists m \in \mathbb{Z}, \ b = am.$$

The expression on the right is a sentence in two variables because the variable $m$ is quantified.

The sentence 'the integer $n$ is not prime' is equivalent to 'there exists a proper divisor of $n$'. Accordingly, we deal with it as follows

$$\exists d \in \mathbb{N}, \; (d|n) \wedge (d \neq 1) \wedge (d \neq n)$$

where we have been careful in excluding trivial divisors.

Let us analyse a famous existence theorem.

**Theorem** (Lagrange[3] 1770). *Every natural number can be written as the sum of four integer squares.*

We consider two instances of this result

$$5 = 2^2 + 1^2 + 0^2 + 0^2 \qquad 7 = 2^2 + 1^2 + 1^2 + 1^2.$$

From the first example, it's clear that some integers are the sum of fewer than four squares. On the other hand, one verifies that 7 does require four squares. We may express these identities without revealing the solution:

$$\exists a, b, c, d \in \mathbb{Z}, \; 7 = a^2 + b^2 + c^2 + d^2.$$

Now, we replace 5 or 7 with an unspecified natural number $n$, and obtain a predicate $\mathscr{L}$ over $\mathbb{N}$, which reads

$$\mathscr{L}(n) := \big(\exists a, b, c, d \in \mathbb{Z}, \; n = a^2 + b^2 + c^2 + d^2\big) \tag{3.22}$$

or

> $n$ is a sum of four integer squares.

To state Lagrange's theorem with symbols, we quantify the remaining variable $n$

$$\forall n \in \mathbb{N}, \; \exists a, b, c, d \in \mathbb{Z}, \; n = a^2 + b^2 + c^2 + d^2. \tag{3.23}$$

Five quantifiers are needed, to turn the predicate $n = a^2 + b^2 + c^2 + d^2$ (a function of five variables) into a boolean constant. Lagrange's theorem can now be translated into a set identity

$$\mathbb{N} = \{n \in \mathbb{N} : \mathscr{L}(n)\}.$$

Indeed *any* theorem on natural numbers has this form, for an appropriate predicate $\mathscr{L}$. (Think about it.)

This analysis clarifies the logical structure of Lagrange's theorem. It also shows that this theorem is better formulated with words than with symbols!

---

[3] Joseph-Louis Lagrange, Italian mathematician and astronomer (1736–1813).

## 3.5   Negating logical expressions

Consider the following sentences

> *Not all primes are odd.*
> *A square integral matrix is not necessarily invertible.*
> *Not all continuous real functions are differentiable.*

These are disguised existence statements

> *There is an even prime.*
> *There is a square integral matrix which is not invertible*
> *There is a continuous real function which is not differentiable.*

Statements of this kind result from the negation of boolean expressions containing quantifiers, to which we now turn.

If $\mathscr{L}$ is a logical expression, its negation is $\neg\mathscr{L}$, which is false if $\mathscr{L}$ is true, and vice-versa. If $\mathscr{L}$ is an expression involving quantifiers and boolean operators, then the expression $\neg\mathscr{L}$ unfolds into an equivalent expression, which we wish to determine.

Theorem 3.1.2 gives us the negation formulae for the main compound expressions

$$
\begin{aligned}
\neg(P \wedge Q) &\Leftrightarrow \neg P \vee \neg Q \\
\neg(P \vee Q) &\Leftrightarrow \neg P \wedge \neg Q \\
\neg(P \Rightarrow Q) &\Leftrightarrow P \wedge (\neg Q).
\end{aligned}
\tag{3.24}
$$

The first two formulae are items (*i*) and (*ii*) of the theorem (de Morgan's laws); the third follows immediately from (*iii*) and (*i*). We see that the negation of an implication does not have the form of an implication.

How does one negate an expression involving quantifiers? We begin with a theorem dealing with the case of a single quantifier; the general case will follow from it.

**Theorem.** *Let $A$ be a non-empty set, let $\mathscr{P}$ be a predicate on $A$, and let*

$$
\mathscr{L} := \forall x \in A, \ \mathscr{P}(x) \qquad \mathscr{M} := \exists x \in A, \ \mathscr{P}(x).
$$

*Then*

$$
\neg\mathscr{L} \Leftrightarrow \exists x \in A, \ \neg\mathscr{P}(x) \qquad \neg\mathscr{M} \Leftrightarrow \forall x \in A, \ \neg\mathscr{P}(x).
\tag{3.25}
$$

PROOF. From (3.18), we have the boolean equivalence

$$\mathscr{L} \Leftrightarrow \mathscr{P}(A) = \{\text{T}\}$$

where $\mathscr{P}(A)$ is the image of $A$ under $\mathscr{P}$ (see section 2.2). But then

$$\neg\mathscr{L} \Leftrightarrow \mathscr{P}(A) \neq \{\text{T}\} \Leftrightarrow \neg\mathscr{P}(A) \neq \{\text{F}\}.$$

Using again (3.18), we obtain

$$\neg\mathscr{L} \Leftrightarrow \exists x \in A, \ \neg\mathscr{P}(x)$$

as claimed.

The second formula is proved similarly. From $\mathscr{M} \Leftrightarrow \mathscr{P}(A) \neq \{F\}$, we obtain

$$\neg\mathscr{M} \Leftrightarrow \mathscr{P}(A) = \{\text{F}\} \Leftrightarrow \neg\mathscr{P}(A) = \{\text{T}\}$$

and hence

$$\neg\mathscr{M} \Leftrightarrow \forall x \in A, \ \neg\mathscr{P}(x).$$

□

It's now easy to deduce the rule for negating expressions with two or more quantifiers. Let $\mathscr{P}$ be a predicate on $A \times B$, and let us consider the expression

$$\mathscr{L} := \forall x \in A, \ \exists y \in B, \ \mathscr{P}(x,y).$$

Then, repeated applications of theorem 3.5 give

$$\begin{aligned}
\neg\mathscr{L} &\Leftrightarrow \neg\big[\forall x \in A, \ (\exists y \in B, \ \mathscr{P}(x,y))\big] \\
&\Leftrightarrow \exists x \in A, \ \neg(\exists y \in B, \ \mathscr{P}(x,y)) \\
&\Leftrightarrow \exists x \in A, \ \forall y \in B, \ \neg\mathscr{P}(x,y).
\end{aligned}$$

It should be clear how similar formulae may be derived, involving any combination of two quantifiers.

Finally, consider the case of several leading quantifiers, e.g.,

$$\mathscr{L} = \forall x_1 \in X_1, \ \exists x_2 \in X_2, \ \ldots \ , \exists x_n \in X_n, \ \mathscr{P}(x_1,\ldots,x_n)$$

where the $X_i$ are sets, and $\mathscr{P}$ is a predicate on the cartesian product $X_1 \times X_2 \times \cdots \times X_n$. Parentheses make it clear that this is a nested array of predicates

$$\mathscr{L} = \forall x_1 \in X_1, \ [\exists x_2 \in X_2, \ [\ldots \ , [\exists x_n \in X_n, \ \mathscr{P}(x_1,\ldots,x_n)]]].$$

Using repeatedly theorem 3.5, from the outside to the inside, we obtain

$$\neg \mathscr{L} \Leftrightarrow \exists x_1 \in X_1, \ \forall x_2 \in X_2, \ \dots \ , \forall x_n \in X_n, \ \neg \mathscr{P}(x_1, \dots, x_n).$$

Namely, the negation of $\mathscr{L}$ is obtained by replacing each $\forall$ with $\exists$, and vice-versa, *without changing their order,* and then replacing $\mathscr{P}$ with $\neg \mathscr{P}$.

   A formal proof of this statement requires the principle of induction, and we postpone it until chapter 7.

EXAMPLE. Let us consider the (false) implication

   *If an integer divides the product of two integers, then it divides one of the factors.*

We negate this statement, in slow motion. We first translate it in symbols

$$\forall n, a, b \in \mathbb{Z}, \ (n|ab) \Rightarrow (n|a \vee n|b).$$

Then we negate it, using theorem 3.5

$$\exists n, a, b \in \mathbb{Z}, \ (n|ab) \wedge (n \nmid a \wedge n \nmid b).$$

Next we express it with a mixture of words and symbols

   *There are integers $n, a, b$ such that $n$ divides the product $ab$, but it does not divide $a$ nor $b$.*

Finally, we dispose of all symbols altogether.

   *There is a natural number which divides the product of two integers without dividing any of the factors.*

EXAMPLE. Lagrange's theorem, expressed symbolically in equation (3.23), would be false if four squares were replaced by three squares. This fact is written symbolically as

$$\exists n \in \mathbb{N}, \ \forall a, b, c \in \mathbb{Z}, \ n \neq a^2 + b^2 + c^2.$$

In words:

   *There is a natural number which cannot be written as the sum of three squares.*

## Exercises

**Exercise 3.1.** Compute the value of the following boolean expressions

1. $(2^{20} > 10!) \lor (2^{10} > 10^3)$
2. $\left(\sqrt{2} > \dfrac{7}{5}\right) \land \left(\sqrt{2} < \dfrac{17}{12}\right)$
3. $1 < 14\,(\sqrt{3} - 2)^2$
4. '47 is the sum of two squares'.

Use only integer arithmetic. In part 4, give only the minimum information needed to establish the value of the expression.

**Exercise 3.2.**

1. Complete the proof of theorem 3.1.2. Hence define the operators $\lor, \Rightarrow, \Leftrightarrow$ in terms of $\neg$ and $\land$.

2. Complete the proof of theorem 3.2.

3. Prove the boolean identity $[(x \land y) \Rightarrow z] \Leftrightarrow [x \Rightarrow (\neg y \lor z)]$.

4. Formulate the negation formulae (3.24) in terms of predicates, and then prove them.
   [*Use theorem 3.2, page 57.*]

**Exercise 3.3.** Consider the boolean expression

$$\forall x \in X, \ \mathscr{P}(x) \Rightarrow \mathscr{Q}(x)$$

where $X$ is a set, and $\mathscr{P}$ and $\mathscr{Q}$ are predicates over $X$. Show that there is a subset $A$ of $X$, such that the expression above is equivalent to

$$\forall x \in A, \ \mathscr{Q}(x).$$

Show that $A$ is a proper subset of $X$ if and only if $\mathscr{P}$ is not constant.

**Exercise 3.4.** Write each statement with symbols, using the quantifier $\exists$.

1. $z \in X + Y$.

2. The integer $n$ is a square.

3. The fraction $a/b$ is not reduced.

4. The unit circle has a rational point.

5. The sets $X$ and $Y$ are not disjoint.

6. The set $X$ is not a subset of the set $Y$.

7. The integer $n$ is not divisible by 3.

**Exercise 3.5.** Write each statement with symbols, using at least one quantifier.

1. The equation $f(x) = 0$ has a rational solution.

2. The equation $f(x) = 0$ has no integer solution.

3. There is no smallest positive rational number.
   [*The set of positive rational numbers is denoted by* $\mathbb{Q}^+$.]

4. There is a rational between any two distinct real numbers.
   [*You must take into account the ordering of such real numbers.*]

5. The function $f : A \to B$ is not injective.

6. The function $f : A \to B$ is not surjective.

7. The function $f : A \to B$ is constant.

8. The function $f : A \to B$ is not constant.

9. Every cubic equation with integer coefficients has a real solution.
   [*Make sure your equation is actually cubic.*]

10. There are integers that can be written in two different ways as the sum of two cubes.[4]

11. The equation $f(x) = 0$ has infinitely many real solutions.

---

[4]This statement refers to the leftmost identity in (3.4), page 50, due to the Indian mathematician Srinivasa Ramanujan (1887–1920).

**Exercise 3.6.** (W. Hodges.) We wish to to build up a set of predicates to describe family relations. You are given the two predicates

'*x is a son of y*'       '*x is a daughter of y*'.

Your task is to write definitions of the following predicates, in some appropriate order so that the later definitions use only the given predicates and earlier definitions. (The order below is just alphabetical.)

> *x is an aunt of y*
> *x is a brother of y*
> *x is a child of y*
> *x is the father of y*
> *x is female*
> *x is a grandchild of y*
> *x is a half-brother of y*
> *x is male*
> *x is the mother of y*
> *x is a parent of y*
> *x is a sister of y*

Use the symbols $x, y, z$, etc., to denote people, words for everything else, and parentheses to specify the order of evaluation of logical operators.

EXAMPLE. To define nephew we must first define of brother and sister.

> *x is a nephew of y* := *There is z such that x is a son of z and*
> *(z is a brother of y or z is a sister of y).*

**Exercise 3.7.** Let $x$ and $y$ be natural numbers, and let $\mathscr{P}(x,y)$ mean: '$x$ is a proper divisor of $y$' (i.e., $x|y$ and $x \neq 1, y$). Thus $\mathscr{P}$ is a predicate over $\mathbb{N} \times \mathbb{N}$.

(a) Express each statement with words, and hence decide if the statement is true or false. [✘, 10]

> 1. $\exists x \in \mathbb{N}, \ \mathscr{P}(x,5)$
> 2. $\exists x \in \mathbb{N}, \ \mathscr{P}(5,x)$
> 3. $\forall x \in \mathbb{N}, \ \mathscr{P}(x,x)$
> 4. $\exists x \in \mathbb{N}, \ \forall y \in \mathbb{N}, \ \mathscr{P}(x,y)$

5. $\forall x \in \mathbb{N} \setminus \{1\}$, $\exists y \in \mathbb{N}$, $\mathscr{P}(x,y)$.

(b) Describe each predicate $\mathscr{F}$ with words, apart from the predicate's variable. [∌, 10]

    1. $\mathscr{F}(y) = \exists x \in \mathbb{N}$, $\mathscr{P}(x,y)$

    2. $\mathscr{F}(y) = \forall x \in \mathbb{N}$, $\neg \mathscr{P}(x,y)$

    3. $\mathscr{F}(x) = \exists y \in \mathbb{N}$, $\mathscr{P}(x,y)$.

(c) Define each set with words. [∌, 10]
[*Consider the previous problem.*]

    1. $\{y \in \mathbb{N} : \exists x \in \mathbb{N}$, $\mathscr{P}(x,y)\}$

    2. $\{y \in \mathbb{N} : \forall x \in \mathbb{N}$, $\neg \mathscr{P}(x,y)\}$

    3. $\{x \in \mathbb{N} : \exists y \in \mathbb{N}$, $\mathscr{P}(x,y)\}$.

**Exercise 3.8.** Let $\mathscr{F}_{\mathbb{R}}$ denote the set of all real functions $\mathbb{R} \to \mathbb{R}$. Turn symbols into words. [∌]
[*Consider theorem 3.1.2, page 55.*]

1. $\{f \in \mathscr{F}_{\mathbb{R}} : f(0) = 0\}$

2. $\{f \in \mathscr{F}_{\mathbb{R}} : \forall x \in \mathbb{R}, f(x) > 0\}$

3. $\{f \in \mathscr{F}_{\mathbb{R}} : \forall x \in \mathbb{Z}, f(x) = 0\}$

4. $\{f \in \mathscr{F}_{\mathbb{R}} : \forall x \in \mathbb{R}, f(f(x)) = x\}$

5. $\{f \in \mathscr{F}_{\mathbb{R}} : \forall x \in \mathbb{R}, x < 0 \Rightarrow f(x) = 0\}$

6. $\{f \in \mathscr{F}_{\mathbb{R}} : \forall x \in \mathbb{R}, 0 \leqslant x \leqslant 1 \Rightarrow f(x) \neq 0\}$

7. $\{f \in \mathscr{F}_{\mathbb{R}} : \forall x \in \mathbb{Z}, (f(x) = 0) \Rightarrow (x = 0)\}$

8. $\{f \in \mathscr{F}_{\mathbb{R}} : \forall x \in \mathbb{Z}, f(x) \neq 0 \Rightarrow x \neq 0\}$

9. $\{f \in \mathscr{F}_{\mathbb{R}} : \forall x \in \mathbb{Q}, f(x) \neq 0 \Leftarrow x \neq 0\}$

**Exercise 3.9.** The following expressions define sets. Turn symbols into words. [∉]
[*Cf. example 3.5.*]

1. $\{n \in \mathbb{Z} : \forall a, b \in \mathbb{Z}, \ (n|ab) \Rightarrow (n|a \vee n|b)\}$

2. $\{a \in \mathbb{Z} : \forall n, b \in \mathbb{Z}, \ (n|ab) \Rightarrow (n|a \vee n|b)\}$

**Exercise 3.10.** In each of the following set definitions, turn words into symbols. [Use the standard definition of sets, and —where appropriate— the quantifier $\forall$. By a *real function* we mean a function from $\mathbb{R}$ to $\mathbb{R}$.]

1. The set of real functions that do not vanish.

2. The set of real functions that vanish at the origin.

3. The set of real functions that do not vanish away from the origin.

4. The set of real functions that vanish identically outside the open unit interval.

5. The set of real functions that vanish at all even integers.

6. The set of real functions that may vanish only at the rationals.
   [*Think about the meaning of 'may vanish'.*]

7. The set of two-by-two invertible real matrices.

8. The set of real matrices representing planar rotations.

**Exercise 3.11.** The following statements are (equivalent to) implications, which may be true or false. For each implication [∉]

   *i)*   state the contrapositive;
  *ii)*   state the converse, and decide whether it is true or false;
 *iii)*   state the negation, and decide whether it is true or false.

1. An integer is also a rational number.

2. The product of two integers is an integer.

3. The square of an even integer is an even integer.

4. If a prime divides the product of two integers, it divides one of them.

5. A cubic equation with real coefficients has at least one real root.

6. Every bounded set of rational numbers is finite.

**Exercise 3.12.** Consider the logical expressions:

1. $\forall n \in \mathbb{N}, \ 1/n \notin \mathbb{N}$

2. $\forall x, y \in \mathbb{R}, \ xy = yx$

3. $\forall n \in \mathbb{N}, \ \sqrt{n} \in \mathbb{R} \smallsetminus \mathbb{Q}$

4. $\forall n \in \mathbb{Z}, \ 2 \mid n(n+1)$

5. $\forall n, m \in \mathbb{Z}, \ (2 \nmid n \wedge 2 \nmid m) \Rightarrow 2 \mid (m+n)$

6. $\forall n \in \mathbb{N}, \ \exists r \in \mathbb{Q}, \ \sqrt{n} < r < \sqrt{n+1}$

7. $\forall x, y \in \mathbb{R}, \ (x < y) \Rightarrow (x^2 < y^2)$

8. $\forall x \in \mathbb{R}, \ \forall y \in \mathbb{R}, \ \exists n \in \mathbb{N}, \ (x > 1) \Rightarrow (x^n > y)$

If an expression is true, state it concisely with words; if it is false, state its negation, first with symbols (using theorem 3.5), then with words. **[⪊]**
[*In part 8, the important quantity is x.*]

**Exercise 3.13.** I have two finite sets, and a function between them. I can only compute the value of the function at each point of its domain, and count and compare sets elements. I need detailed, explicit instructions for answering each of the following questions. **[⪊, 50]**

1. How do I check if my function is surjective?

2. How do I check if my function is injective?

3. Suppose my function is the characteristic function of some set. How do I determine such a set?

**Exercise 3.14.** Consider the statement

*The equation of the unit circle is a relation on the set of real numbers.*

Explain it to someone who is just beginning to learn higher mathematics.

**Exercise 3.15.** Define some interesting predicates on the power set of $\mathbb{Z}$. Do the same for the power set of $\mathbb{Z}[x]$.

# Chapter 4

# Describing functions

In section 2.2 we introduced the basic function terminology. Our vocabulary to describe a function's properties, is still limited to few words (injectivity, surjectivity, invertibility), and we need to expand it. In this chapter we consider attributes of **real functions** $f : \mathbb{R} \to \mathbb{R}$, introducing important terms (such as boundedness and continuity) which are found in more general situations. Then we export this terminology to functions defined over $\mathbb{N}$ —the **real sequences.**

## 4.1  Ordering properties

The real line $\mathbb{R}$ is **ordered**, meaning that for any pair $(x, y)$ of real numbers, precisely one of the three relational expressions

$$x < y \qquad x = y \qquad x > y$$

is true, and the other two are false (see section 3.3.2). We begin to deal with properties that are formulated in terms of ordering.

The simplest attribute of a function concerns the sign of the values it assumes

$$\forall x \in \mathbb{R}, \quad f(x) > 0 \qquad \textbf{positive} \tag{4.1}$$
$$\forall x \in \mathbb{R}, \quad f(x) < 0 \qquad \textbf{negative}. \tag{4.2}$$

If (4.1) is formulated with the non-strict inequality $f(x) \geqslant 0$, then we say that the function is **non-negative**. For the inequality $f(x) \leqslant 0$, the term 'non-positive' is uncommon, and one would normally use **negative or zero**. For example, the exponential function is positive, the absolute value function is non-negative, and the

sine function is neither positive nor negative. Because inequalities are reversed un-
der sign change, if a function $f$ has any of the stated properties, then $-f$ has the
complementary property (e.g, if $f$ is positive, then $-f$ is negative).

We remark that the terms 'non-negative' and 'not negative' have different mean-
ing, the latter being the logical negation of negative. This distinction is very clear
in the symbolic definitions

$$\forall x \in \mathbb{R},\ f(x) \geqslant 0 \qquad \textbf{non-negative}$$
$$\exists x \in \mathbb{R},\ f(x) \geqslant 0 \qquad \textbf{not negative.}$$

(The second symbolic expression is obtained by negating (4.2) according to the
prescription of theorem 3.5.) So a non-negative function is also not negative, but
not vice-versa. Having two similar expressions with rather different meaning can
easily lead to confusion, and one must remain vigilant.

Next we consider how the action of function affects the ordering of the real
line; the order may be preserved, reversed, or a bit of both. There are two com-
peting terminologies, labelled *I* and *II* in the table below. Each has advantages and
disadvantages.

|  | I | II |
|---|---|---|
| $\forall x,y \in \mathbb{R},\ x > y \Rightarrow f(x) > f(y)$ | **increasing** | **strictly increasing** |
| $\forall x,y \in \mathbb{R},\ x > y \Rightarrow f(x) \geqslant f(y)$ | **non-decreasing** | **increasing** |
| $\forall x,y \in \mathbb{R},\ x > y \Rightarrow f(x) < f(y)$ | **decreasing** | **strictly decreasing** |
| $\forall x,y \in \mathbb{R},\ x > y \Rightarrow f(x) \leqslant f(y)$ | **non-increasing** | **decreasing** |

A function that is either increasing or decreasing (strictly or otherwise) is said
to be **monotonic.**

Thus the arc-tangent function is increasing for *I*, and strictly increasing for *II*.
According to terminology *I*, no function can be both increasing and decreasing,
and there are function that are neither, for instance a constant, or the sine function.
The constant functions are the function which are both non-decreasing and non-
increasing.

The disadvantage of *I* is that, as we did above for 'non-negative', we must dif-
ferentiate between 'non-increasing' and the logical negation of increasing ($\exists x,y \in
\mathbb{R},\ (x > y) \wedge (f(x) \leqslant f(y))$).

Terminology *II* eliminates the annoying distinction between the prefixes 'non-'
and 'not', but it introduces a new problem. Now a constant function is both increas-
ing and decreasing, which clashes with common usage. (One wouldn't say 'the
money in my bank account is increasing', if the balance remains the same.)

## 4.2 Symmetries

Real functions may have symmetries, expressing invariance with respect to changes of the independent variable. A function $f$ is **even** and **odd**, respectively, if

$$\forall x \in \mathbb{R}, \ f(-x) = f(x) \qquad \text{and} \qquad \forall x \in \mathbb{R}, \ f(-x) = -f(x),$$

respectively. So the cosine is even, the sine is odd, the exponential is neither even nor odd, and the zero function is both even and odd. The property of being even or odd has a geometrical meaning: graphs of even functions have a mirror symmetry with respect to the ordinate axis, while those of odd functions are symmetrical with respect to the origin. For instance, the function displayed in figure 4.2 is odd. There



Figure 4.1: An odd function

is an easy way of constructing even/odd functions. For any real function $g$, the function $x \mapsto g(x) + g(-x)$ is even, and so is the function $g(f(x))$ for any even function $f$. So the function $x \mapsto g(|x|)$ is even. To construct odd functions, we first define the **sign function**

$$\mathrm{sign}(x) = \begin{cases} +1 & \text{if } x > 0 \\ 0 & \text{if } x = 0 \\ -1 & \text{if } x < 0 \end{cases} \qquad x \in \mathbb{R}. \tag{4.3}$$

The sign function is odd: $\mathrm{sign}(-x) = -\mathrm{sign}(x)$. Then, for any real function $g$, the function $x \mapsto \mathrm{sign}(x) g(|x|)$ is odd, as easily verified. This construct ensures that our function vanishes at the origin, which is a property of all odd functions.

Next we turn to translational symmetry, which is called **periodicity**. A function $f$ is **periodic with period** $T$ if

$$\forall x \in \mathbb{R}, \;\; f(x+T) = f(x) \tag{4.4}$$

for some non-zero real number $T$ (see figure 4.2).



Figure 4.2: A periodic function

For instance, the sine function is periodic with period $T = 2\pi$. If a function is periodic with period $T$, then it is also periodic with period $2T$, $3T$, etc. For this reason, one normally requires the period $T$ to be the smallest positive real number for which (4.4) is satisfied. To emphasise this point, we use the terms **minimal** or **fundamental** period.

If we say that a function $f$ is **periodic** —without reference to a specific period— then the existence of the period must be required explicitly

$$\exists T \in \mathbb{R}^{+}, \;\; \forall x \in \mathbb{R}, \;\; f(x+T) = f(x) \tag{4.5}$$

where the set $\mathbb{R}^{+}$ of positive real numbers was defined in (2.9). Note that the period $T$ must be non-zero (lest this definition says nothing), and **without loss of generality**[1] we shall require the period to be positive. (The case of negative period is dealt with by letting $x = x - T$ in (4.4).)

---

[1] This expression indicates that an inessential restriction or simplification is being introduced. See section 6.8 for another example.

The presence of symmetries reduces the amount of information needed to specify a function. Indeed, if a function is even or odd, knowledge of the function for non-negative values of the argument suffices to characterise it completely. Likewise, if the behaviour of a periodic function is known over any interval of length equal to the period, then the function is specified completely.

It is important to realise that symmetry properties are special, and a function chosen 'at random' will have no symmetry.

## 4.3  Boundedness

A set $X \subset \mathbb{R}$ is **bounded** if there is an interval containing it, namely if[2]

$$\exists a, b \in \mathbb{R}, \ \forall x \in X, \ a < x < b \tag{4.6}$$

or, equivalently, if

$$\exists b \in \mathbb{R}, \ \forall x \in X, \ |x| < b. \tag{4.7}$$

The two definitions imply each other (see exercises).

A function $f$ is **bounded** if its image $f(\mathbb{R})$ is a bounded set. In symbols

$$\exists b \in \mathbb{R}, \ \forall x \in \mathbb{R}, \ |f(x)| < b.$$

For example, the sine function is bounded, while the exponential is not. The periodic function displayed in figure 4.2 is bounded.

A function $f$ is said to be **bounded away from zero** if its reciprocal is bounded. This means that there exists a positive constant $c$ such that $|f(x)| > c$ for all values of $x$. In symbols

$$\exists c \in \mathbb{R}^{+}, \ \forall x \in \mathbb{R}, \ |f(x)| > c.$$

Thus the hyperbolic cosine is bounded away from zero (what could be a value of $c$ in this case?), but the exponential function is not.

## 4.4  Neighbourhoods

A **neighbourhood** of a point $x \in \mathbb{R}$ is any *open* interval containing $x$. The requirement that the interval be open prevents $x$ from being one of the end-points.

---

[2] The term 'interval' is intended in the proper sense —rays are excluded, cf. (2.10.

The neighbourhood concept characterises 'proximity' in a concise manner, that bypasses quantitative statements altogether. A skillful use of this term leads to terse and incisive statements. For instance, the statement

>   *The function $f$ is bounded in a neighbourhood of $x_0$*

means that there exists an open interval containing $x_0$ whose image under $f$ is a bounded set. If we write this statement in symbols

$$\exists a, b \in \mathbb{R}, \ \exists c \in \mathbb{R}^+, \ (a < x_0 < b) \wedge (\forall x \in (a, b), \ |f(x)| < c)$$

we realise just how much information is packed within it.

The statement

>   *Every neighbourhood of $x_0$ contains a point of the set $A$*

is more powerful than what may appear at first sight. The expression 'every neighbourhood of $x_0$' identifies at once infinitely many open intervals, and the ones we are interested in here are the arbitrarily small ones. A more informative sentence would make this explicit

>   *Every neighbourhood of $x_0$, no matter how small, contains a point of the set $A$.*

If $x_0 \in A$, then this statement gives no information, because $x_0$ belongs to every neighbourhood of $x_0$, by definition. The situation is very different if $x_0 \notin A$. For instance, if $x_0 = 0$ and

$$A = \{1, 1/2, 1/3, 1/4, \ldots\} \qquad \text{or} \qquad A = \{\sin(1), \sin(2), \sin(3), \ldots\}$$

then $A$ does not contain zero. It then follows that $A$ contains elements that are **arbitrarily close to zero.** Equivalently, every neighbourhood of zero contains *infinitely many points* of $A$. (Think about it.)

By a **neighbourhood of infinity** we mean a set of the type $\{x \in \mathbb{R} : x > a\}$, where $a$ is a real number. A neighbourhood of $-\infty$ is defined similarly, and indeed the two points at infinity $\pm\infty$ may be identified with the construct $\{x \in \mathbb{R} : |x| > a\}$. For instance, if $A \subset \mathbb{R}$, the expression

>   *There are points of $A$ in any neighbourhood of infinity*

says that the set $A$ is unbounded. The statement

*The function f is constant for all sufficiently large x*

written symbolically as

$$\exists a \in \mathbb{R}, \ \forall x \in \mathbb{R}, \ (x > a) \Rightarrow (f(x) = f(a))$$

says that $f$ is constant in an unspecified neighbourhood of infinity.

The function $f$ has a **maximum** at $x$ if the value of $f$ at $x$ is greater than the value at all other points, namely if

$$\forall y \in \mathbb{R}, \ f(x) > f(y). \tag{4.8}$$

The concept of **minimum** defined similarly. If the strict inequality in definition (4.8) is made non-strict, then we obtain a variant of this concept. The function $f$ has a **local maximum** at $x$ if the property (4.8) holds in some neighbourhood of $x$, rather than in the whole domain. Thus we may write

*The arc-tangent function has no maximum or minimum*
*The sine function has infinitely many local maxima and minima.*

Let $x$ be a real number, and let $f$ be defined in some neighbourhood of $x$, excluding, possibly, the point $x$ itself. If $f$ is unbounded in every neighbourhood of $x$, then we say that $f$ is **singular at $x$**, and the point $x$ is called a **singularity** (or a **singular point**) of the function.

Let us write some sentences about singularities of functions

*A polynomial function has no singularities; the same applies to the sine,*
*the cosine, and the exponential.*
*A rational function has finitely many singularities, which are the roots*
*of the polynomial at denominator.*
*The tangent, secant, and co-secant functions have infinitely many singularities; they form an array of evenly spaced points on the real line.*

Let $\mathcal{N}_x$ be the set of all neighbourhoods of $x$. This is an **abstract set**, whose elements are the open intervals containing $x$. As we did in section 2.1.3, we look for a concrete **representation** for $\mathcal{N}_x$. We identify the elements of this set with the ordered pairs $(a,b)$ of real numbers, specifying the left and right end-points of an interval, respectively[3]. Collecting all these pairs together, we obtain a subset $N_x$ of $\mathbb{R}^2$, given by

$$N_x = \{(a,b) \in \mathbb{R}^2 : a < x < b\} \sim \mathcal{N}_x. \tag{4.9}$$

---

[3]By a Freudian coincidence, the ordered pair $(a,b)$ here represents the open interval $(a,b)$!

The set $N_x$ is the (open) north-west quadrant of the plane, taken with respect to the origin $(x,x)$. In (4.9), we used the symbol '$\sim$' instead of '=', because the sets $\mathcal{N}_x$ and $N_x$ are **equivalent** rather than **equal**, meaning that there is a bi-jective function from $\mathcal{N}_x$ to $N_x$. (Indeed $\mathcal{N}_x$ is *not* a subset of $\mathbb{R}^2$, being a set of open intervals!)

## 4.5   Continuity

Continuity is a fundamental concept in the theory of functions. Loosely speaking, a function is continuous if it has 'no jumps', but this naive definition is only adequate for simple situations. We begin by considering continuity at a specific point of the domain of a function. An informal —yet accurate— characterisation of continuity, is to say that a function is continuous at a point if its value there can be inferred unequivocally from the values at neighbouring points. For example, consider the sign function defined in equation (4.3). Its value at the origin cannot be inferred from the surrounding environment, and even if we defined, say, $\text{sign}(0) = 1$, the ambiguity would remain. To see how bad things can get, consider the real function

$$x \mapsto \begin{cases} 1 & \text{if } x \in \mathbb{Q} \\ 0 & \text{if } x \notin \mathbb{Q}. \end{cases}$$

Given that in *any* neighbourhood of *any* real number there are both rational and irrational numbers, there is no way of inferring the value of this function at a point by considering how the function behaves in the surrounding region.

We define continuity using neighbourhoods, which we denote by the letters $I, J$. We say that $f$ is **continuous at the point** $a$ if

$$\forall J \in \mathcal{N}_{f(a)}, \ \exists I \in \mathcal{N}_a, \ f(I) \subset J. \tag{4.10}$$

The concise definition (4.10) must be analysed in detail. The set $\mathcal{N}_a$ consists of all intervals containing the point $a$ where continuity is being tested. The set $\mathcal{N}_{f(a)}$ consists of all neighbourhoods of $f(a)$, the value of the function at $a$. An arbitrary interval $J$ containing $f(a)$ is given; $J$ could be anything, but we should be ready to handle an absurdly small interval. Our task now is to find an interval $I$ containing $a$ whose image lies inside $J$. The choice of such an interval will clearly depend on $J$. We don't necessarily know what the set $f(I)$ looks like; it's non-empty, but it could be anything. The function is continuous at $a$ provided that, by choosing $I$ small enough, we can ensure that $f(I)$ remains small enough to fit within $J$.

Continuity may be defined without reference to neighbourhoods. In this case however, we must provide qualitative information. The full notation is considerably more complex than in (4.10)

$$\forall \varepsilon \in \mathbb{R}^+, \ \exists \delta \in \mathbb{R}^+, \ \forall x \in \mathbb{R}, \ |x - a| < \delta \Rightarrow |f(x) - f(a)| < \varepsilon$$

and it remains so even if we use short-hands, and omit all references to the set $\mathbb{R}$

$$\forall \varepsilon > 0, \ \exists \delta > 0, \ \forall x \ |x - a| < \delta \Rightarrow |f(x) - f(a)| < \varepsilon.$$

A function $f$ is **continuous** if it's continuous at all points of its domain. To adapt the condition (4.10) to this definition, we need an additional quantifier. Let $A \subset \mathbb{R}$ be the domain of $f$; we say that $f$ is continuous on $A$ if

$$\forall a \in A, \ \forall J \in \mathscr{N}_{f(a)}, \ \exists I \in \mathscr{N}_a, \ f(I) \subset J. \tag{4.11}$$

A real function $f$ is **differentiable at** $a$ if the limit

$$\lim_{x \to a} F(x) \qquad F(x) := \frac{f(x) - f(a)}{x - a} \qquad x \neq a$$

exists. The function $F$ is called the **incremental ratio** of $f$ at $a$. Note that $F$ is not defined at $x = a$. However, if $f$ is differentiable at $a$, then by defining $F(a) := \lim_{x \to a} F(x)$, the function $F$ becomes *continuous* at $a$. Let's sum it up in a sentence.

> *A real function is differentiable at a point if its incremental ratio is continuous at that point.*

A function is said to be **differentiable** if it is differentiable at all points of its domain. A function that is differentiable **sufficiently often** (all derivatives up to a sufficiently high order exist) is said to be **smooth**. The expression 'sufficiently often' is deliberately vague; its precise meaning will depend on the context.

Many elementary real functions are continuous. These include the polynomials, the sine, the cosine, and the exponential function. These functions are also differentiable infinitely often.

## 4.6 Other properties

A real function of the form $f(x) = ax$, where $a$ is a real number, is said to be **linear**. The term linear originates from 'line', and functions of the type $f(x) = ax + b$, are

sometimes referred to as being 'linear', because their graph is a line (the correct term is **affine**). Likewise, a function given by a polynomial will be characterised, in the first instance, by the degree of the polynomial. So we'll speak of a **quadratic, cubic, quartic** function, etc.

Consider the absolute value function, defined as follows

$$|x| = \begin{cases} x & \text{if } x \geqslant 0 \\ -x & \text{if } x < 0. \end{cases} \tag{4.12}$$

This function is not linear, but it is made of two linear pieces, glued together at the origin. Functions made of linear or affine pieces are said to be **piecewise linear** or **piecewise affine** (see equation (4.3) and figure 4.6). A function is **piecewise defined** if its domain is partitioned into disjoint intervals or rays, with the function being specified independently over each interval. The properties of a piecewise-defined function may fail at the end-points of the intervals of definition. Under such circumstance, terms such as **piecewise increasing**, **piecewise continuous**, **piecewise differentiable, piecewise smooth** may be used —see figure 4.4. A piecewise constant function is called a **step-function**.



Figure 4.3: A continuous piecewise affine function.

EXAMPLE. If $x$ is a real number, we define the **floor** of $x$, denoted by $\lfloor x \rfloor$, to be the largest integer not exceeding $x$. Similarly, the **ceiling** function $\lceil x \rceil$ represents the smallest integer not smaller than $x$. Floor and ceiling functions are prominent examples of step functions. Closely connected to the floor function is the **fractional part** of a real number $x$, denoted by $\{x\}$. (The notational clash with the set having
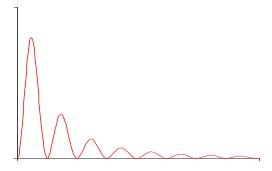
Figure 4.4: Left: a step function. (The vertical segments are just a guide to the eye; they are not part of the graph of the function.) Right: a periodic continuous function which is piecewise differentiable.

$x$ as its only element is one of the most spectacular in mathematics!) The fractional part is defined via the equation

$$x = \lfloor x \rfloor + \{x\}$$

from which it follows that $0 \leqslant \{x\} < 1$. This function is discontinuous and piecewise affine.

EXAMPLE. Describe the following function: [✗]



*This is a smooth function, which is bounded and non-negative. It features an infinite sequence of evenly spaced local maxima, whose height decreases monotonically to zero for large arguments. There is one zero of the function between any two consecutive local maxima.*

We rewrite it, borrowing some terminology from physics.

> *This function displays regular oscillations of constant period, with amplitude decreasing monotonically to zero.*

The following functions have a behaviour that is qualitatively similar to that displayed above

$$f(x) = \sin(x)^2 e^{-x} \qquad\qquad f(x) = \frac{1 + \sin(x)}{1 + x^2}.$$

## 4.7   Describing real sequences

A sequence can be thought of as a function defined over the natural numbers, or, more generally, over a subset of the integers (section 2.4). Using this analogy, some terminology introduced for real functions translates literally to real sequences. So the terms

> positive,   negative,   increasing,   decreasing,   monotonic,
> periodic,   bounded

have the same meaning for sequences as they have for functions. So we write

> *The sequence of primes is positive, increasing and unbounded.*

Other terms required amendments, or are simply not relevant to sequences. For instance, injectivity is not used —we simply say that the terms of the sequence are distinct— while surjectivity is rarely significant (why?). Invertibility is used in a different sense —see section 8.4. The terms 'even' and 'odd' are still applicable to doubly-infinite sequences.

A sequence which settles down to a periodic pattern from a certain point on is said to be **eventually periodic**. More precisely, a sequence $(x_1, x_2, \ldots)$ is eventually periodic if the set $\{x_1, x_2, \ldots\}$, which is the image of $\mathbb{N}$ under the corresponding function, is finite. (Think about it.) If the periodic pattern consists of a single term, then the sequence is said to be **eventually constant**.

The the idea of continuity does not apply to sequences, because the neighbourhood concept loses relevance on the integers. The notable exception is the **point at infinity**. The **neighbourhoods of infinity** are infinite sets of the form $\{n \in \mathbb{N} : n > M\}$, for some integer $M$.

Accordingly, we have the notion of 'continuity at infinity': a real sequence is said to **converge**, if —as a function defined over the integers— it's **continuous at**

**infinity**. Adapting the definition of continuity (4.10) to the present situation, we say that the real sequence $(a_k)$ (which is a function $a : \mathbb{N} \to \mathbb{R}$) converges to the limit $c$, if

$$\forall J \in \mathcal{N}_c, \ \exists I \in \mathcal{N}_\infty, \ a(I) \subset J.$$

Reverting to sequence terminology, this says that given any neighbourhood $J$ of $c$, we have $a_k \in J$ for all sufficiently large $k$. The expression '**for all sufficiently large** $k$' means that there is an integer $M$ such that this property holds for all $k > M$. As for continuity, the symbolic description becomes more complicated without neighbourhoods

$$\forall \varepsilon > 0, \ \exists M, \ \forall k > M, \ |a_k - c| < \varepsilon.$$

## Exercises

**Exercise 4.1.** Consider the following implications, where $f$ is a real function.

1. If $f$ is decreasing, then $-f$ is increasing.

2. If $f$ is decreasing, then $|f|$ is increasing.

3. If $|f|$ is increasing, then $f$ is monotonic.

4. If $f$ is odd, then $f^2$ is even.

5. If $f$ is unbounded, then $f$ is surjective.

6. If $f$ is continuous, then $|f|$ is continuous.

7. If $f$ is differentiable, then $|f|$ is differentiable.

8. If $f$ is periodic, then $f^2$ is periodic.

Of each implication:

    (a) state the converse, and decide whether it's true or false.
    (b) state the contrapositive, and decide whether it's true or false.

**Exercise 4.2.** Prove that the definitions of boundedness (4.6) and (4.7) are equivalent. More precisely, let $\mathscr{B}$ and $\mathscr{B}'$ be the sets of all subsets of $\mathbb{R}$ satisfying (4.6) and (4.7), respectively. Show that $\mathscr{B} = \mathscr{B}'$, that is, that $(\mathscr{B} \subset \mathscr{B}') \wedge (\mathscr{B}' \subset \mathscr{B})$.
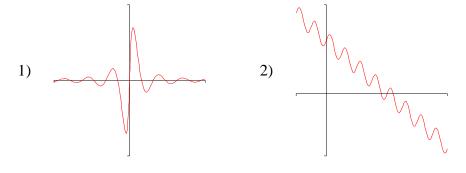
**Exercise 4.3.** Consider the definition (4.5) of a periodic function. What happens if we place the quantifiers in reverse order?
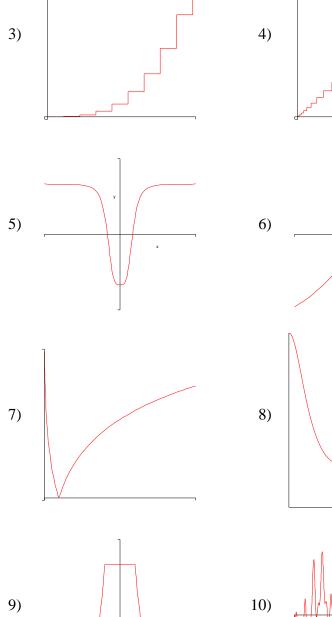
$$\forall x \in \mathbb{R}, \ \exists T \in \mathbb{R}^*, \ f(x+T) = f(x) \tag{4.13}$$
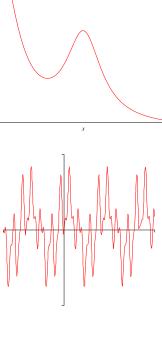$$\forall x \in \mathbb{R}, \ \exists T \in \mathbb{R}^+, \ f(x+T) = f(x). \tag{4.14}$$

(a) Find a function that satisfies (4.13) but not (4.5).

(a) find a function that satisfies (4.14) but not (4.5).

(c) Characterize the set of functions that satisfy (4.13). Do the same for (4.14).

**Exercise 4.4.** Describe the behaviour of each of the following functions.  [♉, 30].

3)

4)

5)

6)

7)

8)

9)

10)

**Exercise 4.5.** In all cases of the previous problem, define a real function $x \mapsto f(x)$ whose behaviour is qualitatively similar to the one displayed.
[*Experiment with Maple.*]

**Exercise 4.6.** Define a real function $f$ with the stated properties.

1. $f$ is increasing, and its derivative is decreasing.

2. $f$ is bounded, and its derivative is unbounded.

3. $f$ is unbounded, and its derivative is bounded.

4. $f$ is discontinuous and its absolute value is continuous.

**Exercise 4.7.** Define a real function $f$ such that its image $f(\mathbb{R})$ is equal to

$$i) \quad [0,1), \qquad\qquad ii) \quad \mathbb{N}, \qquad\qquad iii) \quad \mathbb{Q}.$$

**Exercise 4.8.** Express each of the following statements with symbols.

1. The sequences $(a_k)$ and $(b_k)$ are distinct.

2. The sequence $(a_k)$ is eventually constant.

3. The sequence $(a_k)$ is eventually periodic.

4. The sequence $(a_k)$ has infinitely many negative terms.

5. Each term of the sequence $(a_k)$ appears infinitely often.

# Chapter 5

# Writing effectively

In this chapter we consider some techniques for writing mathematics. We deal with small-scale features: choosing an appropriate terminology and notation, writing clear formulae, mixing words and symbols, writing definitions, presenting a concept. We are not yet concerned with writing a proof, or organising a document.

## 5.1 Choosing words

As we learn how to write mathematics, our first aim is to achieve total accuracy. We analyse some typical mistakes and imprecisions which result from a poor choice of words.

BAD: the equation $x - 3 < 0$

GOOD: the inequality $x - 3 < 0$

BAD: the equation $x^2 - 1 = (x - 1)(x + 1)$

GOOD: the identity $x^2 - 1 = (x - 1)(x + 1)$

BAD: the interval $[0, \infty)$

GOOD: the ray $[0, \infty)$ (the infinite interval $[0, \infty)$)

BAD: the solution of $x^k = x$

GOOD: a solution of $x^k = x$

BAD: the function $f(x)$

GOOD:  the function $f$

 BAD:  the area of the unit circle

GOOD:  the area of the unit disc

 BAD:  the function $g(A)$ of the set $A$

GOOD:  the image $g(A)$ of the set $A$

 BAD:  the absolute value is positive

GOOD:  the absolute value is non-negative

 BAD:  the coordinates of a complex number

GOOD:  the real and imaginary parts of a complex number

 BAD:  the exponential function crosses the vertical axis at a positive point

GOOD:  the graph of the exponential function intersects the ordinate axis at a positive
       point

Once our writing is accurate, we begin to refine the choice of words, to differentiate meaning, or just avoid repetition. There is a rich dictionary at our disposal.

The word **element** denotes a special kind of subsidiary relationship. If $A = \{a_1, \ldots, a_n\}$ is a set, then it is appropriate to write that $a_j$ is an element of $A$. The same applies to sequences, where the word **term** may be used as an alternative to 'element'. The word 'term' is also used to denote the operands in sums or products of the elements of sequences.

A different terminology is used for geometrical objects. If $V = (v_1, \ldots, v_n)$ is a vector, then $v_j$ is a **component** of $V$ (even though a vector is just a finite sequence). If $L$ is a line, and $x \in L$, then $x$ is a **point** of $L$. For any set, the term **point** is a useful alternative to 'element'; you may also use **member**, to avoid repetition.

The term **variable** is used in connection with functions and equations. In the first case, it has the same meaning than **argument**, which refers to the function's input data. In the second case it means **unknown** —a quantity whose value is to be found. Polynomials and rational functions may represent functions or algebraic objects. In the latter case, the term **indeterminate** is preferable to variable.

The term **parameter** is used to identify a variable which is assigned a value that remains fixed in the subsequent discussion. For example, the indefinite integral of a function is a function which depends on a parameter, the integration constant:

$$\int g(x)dx = f(x) + c$$

The two symbols $x$ and $c$ play a very different role here, so we have a **one-parameter family** of functions, rather than a function of two variables.

We now turn to the vocabulary on functions. The word **function** is invariably used for real functions and for real-valued function (the co-domain is $\mathbb{R}$). The term **map** (or **mapping**) is preferable for more general domains and co-domains, and, in particular, in higher dimensions. So a function $f : \mathbb{C} \to \mathbb{C}$ or $f : \mathbb{R}^n \to \mathbb{R}^n$ will be called a mapping. The term **operator** is used to denote an important function of two variables, invariably represented by a symbol rather than a letter ($+, \Rightarrow, \in$, etc.), and short-hand notation ($x + y$, rather than $+(x, y)$). Examples are the **arithmetical operators,** acting on numbers, the **set operators**, acting on sets, and the **logical operators** acting on boolean quantities. The term **binary** means that these operators are functions of two variables. There are also **unary operators,** with a single argument, such as the function that changes the sign of a number, or takes the complement of a set.     The term **operator** is also used to describe functions acting on functions, which produce other functions. The **differentiation operator** is a well-known example. A real-valued function acting on functions is called a **functional**. So definite integration is a functional.

In logic, the terms **predicate, boolean function,** and **characteristic function** represent the same thing. Which term should we use? 'Characteristic function' should be used if there is explicit reference to the associated set, and 'predicate' (or 'boolean function') otherwise. If the input values of a function are also boolean, then there is a strong case for using 'boolean function'. The following examples illustrate the usage of these terms

*The negation of a predicate is a predicate.*

*Let $\chi$ be the characteristic function of the prime numbers.*

*Let us consider the boolean function $(P, Q) \mapsto \neg P \wedge Q$.*

*The predicate '$n \mapsto (2 \mid n)$' is the characteristic function of the even integers.*

*A relational operator is a boolean function of two variables.*

# 5.2   Choosing symbols

Choosing mathematical notation is difficult. Traditionally, mathematicians have been reluctant to accept standardisation of notation, to a degree unknown in other disciplines. Indeed, the ability to adjust quickly to new notation is regarded as

one of the skills of the trade.  The reality is somewhat different: absorbing new notation requires an effort, and most people would be glad to avoid it. So, in order to communicate mathematical ideas without confusing or alienating an audience, the notation must be simple, logical, and consistent.

How does one choose symbols? There are two golden rules:

– DO NOT INTRODUCE UNNECESSARY SYMBOLS.

– DEFINE EACH SYMBOL BEFORE IT'S USED.

Once defined, a symbol should be used consistently: never use the same notation for two different things, or two symbols for the same thing, even if these instances appear far apart in a document. Don't write '$A_j$, for $1 \leqslant j < n$' in one place and '$A_k$, for $1 \leqslant k < n$' in another, unless there is a good reason for doing this. Such small inconsistencies introduce in the text a certain amount of 'notational pollution'. As the pollution piles up, reading becomes tiresome.

I will now offer some guidelines on how to choose symbols. These are not rigid prescriptions; they may be adapted to one's taste, or even rejected altogether. What's important is to develop awareness of notation, and to make conscious decisions about it.

SETS.  Represent sets by capital letters, Roman or Greek, such as

$$S \qquad \mathscr{A} \qquad \Omega.$$

The large variety of fonts available in modern typesetting systems increases our choice. When dealing with generic sets, then $A, B, C$ or $X, Y, Z$ are good symbols. For specific sets, choose a symbol that will remind the reader of the nature of the set. So, for an **alphabet**, $\{a, b, c, \ldots\}$, the symbols $A, \mathscr{A}$ are obvious choices; likewise, $F$ is appropriate for a set of functions, etc.

Lower-case symbols like $x, y$ represent the elements of a set. So $x \in A$ is a good notation, $X \in A$ is bad, and $X \in a$ is very bad. If more than one set is involved, it may be helpful to use matching symbols. Thus

$$a \in A \qquad b \in B \qquad c \in C$$

is more coherent than

$$x \in A \qquad y \in B \qquad z \in C.$$

Some delicacy is required in the case of sets of sets. Consider the expressions

$$f(A \cap B) \qquad\qquad f(x \cap y).$$

The left-most expressions will be interpreted as the image of the intersection of two sets under a function. In this case $A \cap B$ represents a subset of the domain of $f$. However, if the domain of $f$ is a set of sets (e.g., a power set), then this expression becomes dangerously ambiguous. The right-most expression removes this ambiguity. The combination of standard symbols $x, y$ for variables, and a set operator gives a clear signal that the argument of the function is an element, rather than a subset, of the domain.

INTEGERS. The choice of a symbol for an integer begins from the middle region of the Roman alphabet

$$i, j, k, l, m, n \tag{5.1}$$

particularly if an integer is used as subscript or superscript[1]. However, use $p$ for a prime number, and $q$ if there is a prime different from $p$. The list of adjacent letters (5.1) cannot be extended; the preceding symbols, $f, g, h$, are usually reserved for functions' names —see below— while the symbol that follows, $o$, is rarely used, not only for its resemblance to 0 (zero), but also because it has an established meaning in asymptotic analysis, where it appears in expressions of the form $o(\log x)$. Capital letters in the list (5.1) may be used in combination with lower-case letters to denote a range of integers $n = 1, \ldots, N$, or to represent large integers. We'll return to this point in section 5.3, in connection with sums and products.

RATIONALS. For rational numbers, use lower-case Roman letters in the ranges $a$–$e$, or $p$–$z$. The notation

$$r = \frac{m}{n}$$

is good, because '$r$' reminds us of 'rational', while numerator and denominator conform to the convention for integers. If there is more than one rational, use adjacent symbols, $s, t$ in this case.

REAL NUMBERS. For real numbers, use the same part of the Roman alphabet as for rationals, or the Greek alphabet:

$$\alpha, \beta, \gamma, \ldots$$

If there are both rational and real numbers, and if the distinction between them is important, use Roman for the rationals and Greek for the reals.

---

[1]Physicists use Greek letters for subscripts and superscripts.

Some Greek symbols have preferential meaning: small quantities are usually represented by $\varepsilon, \delta$, while for angles one uses $\phi, \varphi, \theta$. Some important real numbers have dedicated symbols:

$$\begin{aligned}
\pi &= 3.141592653\ldots & &\textbf{Archimedes' constant} & &(5.2)\\
e &= 2.718281828\ldots & &\textbf{Napier's constant}\\
\gamma &= 0.333177924\ldots & &\textbf{Euler-Mascheroni constant.}
\end{aligned}$$

COMPLEX NUMBERS. Complex numbers tend to occupy the end of the Roman alphabet, and your first choice should be $z$ or $w$. On the complex plane, we write $z = x + iy$, where $x$ and $y$ are the real and imaginary parts of $z$, and $i$ is **the imaginary unit**. (However, number theorists use $\sqrt{-1}$, not $i$.) In polar coordinates, the standard notation is $z = \rho e^{i\theta}$. Be careful not to use $i$ for any other purpose, such as a summation index.

UNKNOWNS. The quintessential symbol for an equation's unknown is $x$, invariably followed by $y$ and $z$ if there are other unknowns. For large numbers of unknowns, it is necessary to use sequence notation $x_1, \ldots, x_n$. The same notation is also appropriate for the indeterminates of a polynomial, or the arguments of a function.

COMPOSITE OBJECTS. Composite objects (groups, graphs, matrices) are best represented with capital letters, both Roman and Greek. So use $G$ or $\Gamma$ for a group or a graph, and $M$ for a matrix. If you have two groups, use adjacent symbols, like $G$ and $H$. As with sets, for these objects' components consider using matching symbols, e.g., $g \in G$. A notable exception are graphs, where $v$ and $e$ are invariably used for vertices and edges, respectively.

FUNCTIONS. The default choice for a function's name is, of course, $f$; if there is more than one function, use the adjacent symbols $g, h$. These lower-case symbols work well with any number of variables: $f(x), f(x, y, x), f(x_1, \ldots, x_n)$. If the codomain of a function is a cartesian product, then the function is a composite object —a vector, in fact— for which capital letters are appropriate. So, a real function of two variables may be specified as

$$F : \mathbb{R}^2 \to \mathbb{R}^2 \qquad (x, y) \mapsto (f_1(x, y), f_2(x, y)).$$

Greek symbols, either capital or lower-case, are also commonly used for functions' names. The contrast between Roman and Greek symbols may be exploited to separate out the symbols' roles, as in $\mu(x)$ or $f(\lambda)$.

Some famous functions are named after, and represented by, a symbol (often a Greek one), thereby creating a strong bond between object and notation. The best known are **Euler's gamma function** $\Gamma$

$$\Gamma(s) = \int_0^\infty x^{s-1} e^{-x} dx$$

(the extension of the factorial function to complex arguments), and **Riemann's zeta-function** $\zeta$

$$\zeta(s) = \sum_{n=1}^\infty \frac{1}{n^s}. \tag{5.3}$$

There is a peculiar notation for this function: its complex argument $s$ is commonly written as $s = \sigma + i\tau$, with $\sigma$ and $\tau$ real numbers. Other functions with dedicated notation are Euler's $\varphi$-function (see section 5.3), Dedekind's $\eta$-function, Kroneker's $\delta$-function, Weierstrass' $\mathscr{P}$-function, Lambert $\mathscr{W}$-function, etc.

SEQUENCES AND VECTORS. Sequences pose specific notational problems, due to the presence of indices. Consider the various possibilities listed in (2.18), page 29: which one should we choose? Be guided by the principle of economy: a symbol should be introduced only if it's strictly necessary. So the notation $(a_k)$ is quite adequate for a generic sequence, or if the specific properties of the sequence (the initial value of the running index, its finiteness) are not relevant. When more information is needed, the notation $(a_k)_{k \geqslant 1}$ is more economical than $(a_k)_{k=1}^\infty$, but the latter may be a better choice if it is to be contrasted with $(a_k)_{k=1}^n$. In turn, the latter notation is not as friendly as $(a_1, \ldots, a_n)$, although it is more concise.

If a sequence is referred to often, even the stripped down notation $(a_k)$ could become heavy, and it may be advisable to allocate a symbol for the sequence

$$a = (a_1, a_2, \ldots) \qquad \mathbf{v} = (v_1, \ldots, v_n).$$

As usual, we have employed matching symbols, using, respectively, a minimalist lower-case Roman character, and a lower-case boldface character. The latter usually represents a vector. When using ellipses, two or three terms of the sequence usually suffice, but there are circumstances where more terms or a different arrangement of terms are needed.

For example, in the expression

$$(1+x, 1+x^2, \ldots, 1+x^{2^k}, \ldots)$$

the insertion of the general term removes any ambiguity, while the ellipsis on the right suggests that the sequence is infinite —cf. expressions (2.17). The notation

$$(a_1, \ldots, a_{k-1}, a_{k+1}, \ldots, a_n)$$

denotes a subsequence of a finite sequence, obtained by removing the $k$-th term. Note that this arrangement does not allow us to remove the first or the last terms.

Things get complicated with sequences of sequences. This construction is not at all unusual; for instance, we may have a sequence of vectors, whose components need to be referred to explicitly. We may write

$$V = (V_1, V_2, \ldots) \qquad \text{or} \qquad \mathbf{v} = (\mathbf{v}_1, \mathbf{v}_2, \ldots).$$

Let $V_k$ (or $\mathbf{v}_k$) be the general term of our sequence. How are we to represent its components? As usual, we choose the matching symbol $v$, with a subscript indicating the component. However, the integer $k$ has to appear somewhere, and its range must be specified. It is advisable to keep $k$ out of the way it as much as possible

$$V_k = (v_1^{(k)}, \ldots, v_n^{(k)}) \qquad k \geqslant 1.$$

(For a variant, replace the inequality $k \geqslant 1$ with the expression $k = 1, 2, \ldots$.) The parentheses are obviously needed here, for otherwise $v_i^k$ would be interpreted as $v_i$ raised to the $k$-th power. However, it may just happen that we need to raise the vector components to some power. Clearly, we can't use adjacent superscripts $v_3^{(2)4}$, so parentheses are needed, but the straightforward notation $(v_3^{(2)})^4$ is awkward. For a more elegant solution, we represent $k$ as an additional subscript, adopting, in effect, matrix notation

$$V_k = (v_{1,k}, \ldots, v_{n,k}) \qquad k \geqslant 1.$$

We remark that with vectors the multiplication symbols '$\cdot$' and '$\times$' are reserved for the scalar and vector products, respectively —cf. (2.8), page 17. Hence for scalar multiplication we must use juxtaposition

$$a(bV \cdot cW) \qquad\qquad x\mathbf{v} \times y\mathbf{u}.$$

DERIVED SYMBOLS.  Closely related objects require closely related notation. Proximity in the alphabet, e.g., $x, y, z$ may be used for this purpose. For a stronger bond, the meaning of a symbol may be modified using subscripts, superscripts and other attachments

$$A^* \qquad \overline{\eta} \qquad n^+ \qquad \underline{h} \qquad \tilde{e} \qquad \Omega_- \qquad Z_r.$$

The following variants of the symbol $\mathbb{R}$ are commonly used:

$$\mathbb{R}^+ = \{x \in \mathbb{R} : x > 0\} \qquad \mathbb{R}^* = \mathbb{R} \smallsetminus \{0\} \qquad \overline{\mathbb{R}} = \mathbb{R} \cup \{\infty\}.$$

It must be noted that there is no general agreement on the meaning of such adjustments. Thus, for a set, the overbar denotes the so-called **closure** —adjoining to a set all its limit points (the transformation from $\mathbb{R}$ to $\overline{\mathbb{R}}$ is essentially a closure operation). However, for complex numbers the overbar denotes complex conjugation. If $f$ is a function, then $f'$ is the derivative of $f$, but for sets a prime indicates taking the complement.

EXAMPLE. The following sentences illustrate the use of derived symbols.

Let $f : X \to X$ be a function, and let $x^* = f(x^*)$ be a fixed point of $f$.

We consider the endpoints $x_-$ and $x_+$ of an interval containing $x$.

Let $f$ be a real function, and let

$$f^+ : \mathbb{R} \to \mathbb{R} \qquad x \mapsto \begin{cases} f(x) & \text{if } f(x) \geqslant 0 \\ 0 & \text{if } f(x) < 0. \end{cases}$$

EXAMPLE. Consider the polynomial in two indeterminates

$$f(x, z) = -z^2 + xz + 1.$$

For fixed value of $z$, we have a polynomial in $x$. In such a situation, we adopt a notation that emphasises the different role played by the symbols $x$ and $z$. Accordingly, we replace $z$ with $a$, to keep it far apart from $x$ in the alphabet, and then we rewrite the expression above as

$$f_a(x) = ax + 1 - a^2 \qquad a \in \mathbb{R}. \tag{5.4}$$

We now have a **one-parameter family** of linear polynomials in $x$. For fixed $a$, the equation $y = f_a(x)$ is the cartesian equation of a line, and we also have a one-parameter family of lines on the plane. Plotting some of these lines reveals a hidden structure: they form the **envelope** of a parabola. Likewise, if we fix $x = a$ we obtain a one-parameter family of quadratic polynomials $g_a(z) = -z^2 + az + 1$.
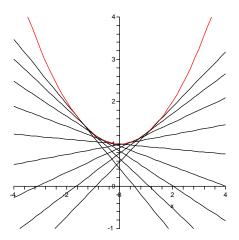
Figure 5.1: The one-parameter family of lines $y = f_a(x)$, where $f_a$ is given by (5.4). These lines are tangent to the parabola $y = x^2/4 + 1$.

## 5.3   The sigma-notation

The notation for sums

$$\sum_{k=1}^{n} a_k = a_1 + a_2 + \cdots + a_n \qquad \sum_{k=1}^{\infty} a_k = a_1 + a_2 + \cdots, \qquad (5.5)$$

was introduced by Fourier[2] (see section 2.4.1). This is called the (delimited) **sigma-notation**, as it makes use of the capital Greek letter with that name. The index of summation is a **dummy variable,** meaning that its identity is irrelevant

$$\sum_{k=1}^{n} k^2 = \sum_{j=1}^{n} j^2 = 1^2 + 2^2 + \cdots + n^2.$$

The summation index is invariably one of the six lower-case roman letters listed in (5.1). The reader is strongly advised to choose a summation symbol, and then stick to it, unless there is a good reason to change symbol.

A **double sum** is defined as follows

$$\sum_{j=1}^{J} \sum_{k=1}^{K} a_{j,k} \stackrel{\text{def}}{=} \sum_{j=1}^{J} \left( \sum_{k=1}^{K} a_{j,k} \right)$$

---

[2]Jean Baptiste Joseph Fourier, French mathematician and physicist (1768–1830).

$$= \sum_{k=1}^{K} a_{1,k} + \sum_{k=1}^{K} a_{2,k} + \cdots + \sum_{k=1}^{K} a_{J,k}.$$

The sum in parentheses is a function of the outer summation index $j$; this sum is performed repeatedly, each time with a different value of $j$. The use of matching symbols for the index and the upper bound of summation is particularly appropriate here. If the two ranges of summations are independent, inner and outer sums can be swapped. The commutative and associative laws of addition ensure that the value of the sum will not change. We illustrate this process with an example

$$\begin{aligned}
\sum_{j=0}^{1} \sum_{k=1}^{3} a_{j,k} &= (a_{0,1} + a_{0,2} + a_{0,3}) + (a_{1,1} + a_{1,2} + a_{1,3}) \\
&= (a_{0,1} + a_{1,1}) + (a_{0,2} + a_{1,2}) + (a_{0,3} + a_{1,3}) \\
&= \sum_{k=1}^{3} \sum_{j=0}^{1} a_{j,k}
\end{aligned}$$

If the indices in a double sum have the same range, then they may be grouped together. The following expressions are equal

$$\sum_{i=1}^{N} \sum_{j=1}^{N} a_{i,j} \qquad \sum_{i,j=1}^{N} a_{i,j}. \tag{5.6}$$

To specify the range of summation, there are variants to the delimited sigma-notation (5.5). If a sum is unrestricted, then range information may be omitted altogether, as in the following sum

$$\sum_{k} \binom{n}{k} = 2^{n} \quad n \geq 0.$$

(This sum has in fact only finitely many non-zero terms.) The summation range may also be specified by inequalities placed below the summation symbols. Expression (5.5) and (5.6) may we rewritten as

$$\sum_{1 \leqslant k \leqslant n} a_{k}, \qquad \sum_{k \geqslant 1} a_{k} \quad \text{and} \quad \sum_{1 \leqslant j,k \leqslant n} a_{j,k}, \tag{5.7}$$

respectively. The benefits of this notation become evident if we need to change summation index. Consider the following manipulation

$$\sum_{-2 \leqslant k \leqslant n-3} 2^{k+2} = \sum_{0 \leqslant k+2 \leqslant n-1} 2^{k+2} = \sum_{0 \leqslant k \leqslant n-1} 2^{k} = 2^{n} - 1. \tag{5.8}$$

After adding 2 to each term in the inequalities, we have simply replaced $k+2$ with $k$, obtaining the sum of a geometric progression. The latter is then evaluated explicitly. With this notation, the change in summation index is effortless.

The greatest generality is achieved by the **standard form** of the sigma-notation

$$\sum_{\mathscr{P}(k)} a_k \tag{5.9}$$

where $\mathscr{P}$ is a predicate over $\mathbb{Z}$. The range of summation consists of those values of $k$ for which $\mathscr{P}(k)$ is TRUE. The unrestricted summation becomes

$$\sum_{k \in \mathbb{Z}} a_k.$$

If we interpret the inequalities in (5.7) as defining a predicate, then we may alter the range of summation by adding conditions

$$\sum_{0 < |k| \leqslant 2} a_k \;=\; a_{-2} + a_{-1} + a_1 + a_2$$

$$\sum_{\substack{1 \leqslant k \leqslant 12 \\ \gcd(k,12)=1}} a_k \;=\; a_1 + a_5 + a_7 + a_{11}$$

$$\sum_{\substack{1 \leqslant k \leqslant 10 \\ k \text{ prime}}} a_k \;=\; a_2 + a_3 + a_5 + a_7.$$

EXAMPLE. For any natural number $n$, we let $\varphi(n)$ be the number of positive integers smaller than $n$ and relatively prime to it (with $\varphi(1) = 1$). Thus $\varphi(12) = 4$. This is **Euler's $\varphi$-function** of number theory, which is defined in symbols as follows

$$\varphi(1) = 1 \qquad \varphi(n) = \sum_{\substack{1 \leqslant k < n \\ \gcd(k,n)=1}} 1, \qquad n > 1.$$

Simply by letting $a_k = 1$ in (5.9), we have turned a summation into a function that counts the elements of the set defined by the predicate $\mathscr{P}$. This concise construction illustrates the effectiveness of the sigma-notation. Alternatively, Euler's function may be defined as the cardinality of a set, using the cardinality symbol '#' (see section 2.1)

$$\varphi(1) = 1 \qquad \varphi(n) = \#\{k \in \mathbb{N} : k < n, \ \gcd(k,n) = 1\} \qquad n > 1.$$

The constructs introduced above for sums are also applicable to products, as well as to combinations of sums and products.

## 5.4  Improving formulae

When the physicist Stephen Hawking was writing his book 'A brief history of time', an editor warned him that for every equation in the book, the readership would be halved. So he included a single equation. This anecdote conveys the hostility that the general public holds towards formulae. Although mathematicians are trained to deal with them, no one likes to struggle with an obscure collection of symbols. In this section we explore some techniques to improve the clarity of a formula, through presentation, notation, and layout.

Formulae are either embedded in the text, or displayed, and in a document one normally finds both arrangements. In either case, a formula must obey standard punctuation rules. The following passage features several embedded formulae, exemplifying an appropriate punctuation.

> *For each $x \in X$, we have the decomposition $x = \xi + \lambda$, with $\xi \in \Xi$ and $\lambda \in \Lambda$; accordingly, we define the function $P : X \to \Xi$, $x \mapsto \xi$, which extracts the first component of $x$.*

The punctuation generates rests as it would in an ordinary English sentence. The two components in the definition of the function $P$ are separated by a comma, which would not be necessary in a displayed formula, see (2.13), and below. To assist the reader, the definition of $P$ is given twice, first with symbols, then with words.

The following example shows the use of punctuation within a displayed formula

$$a_0 = 1; \quad a_{k+1} = \begin{cases} a_k^2 - 1, & 1 \leqslant k < 10; \\ a_k^2, & k \geqslant 10. \end{cases}$$

This is a fully punctuated formula, which some may find pedantic. For a lighter delivery, some punctuation may be replaced by increased spacing, or by words, as follows

$$a_0 = 1 \qquad a_{k+1} = \begin{cases} a_k^2 - 1 & \text{if } 1 \leqslant k < 10 \\ a_k^2 & \text{if } k \geqslant 10. \end{cases}$$

A displayed formula is normally embedded within a sentence, and the punctuation at the end of a formula must be appropriate to the structure of the sentence. In particular, if a sentence terminates at a formula —as in the example above— the full stop at the end of the formula must always be present.

EXAMPLE. The following untidy formula

$$f(x) = \frac{14x - 2x^3 - 2x^2 + 14}{-2x - 4}$$

could represent a typical unprocessed output of a computer-algebra system. It contains redundant information (a common factor between numerator and denominator), the monomials at numerator are not ordered, and there are too many negative signs. The properties of $f(x)$ are not evident from it. There are many ways to improve the layout:

$$f(x) = \frac{x^3 + x^2 - 7x - 7}{x+2}$$

$$f(x) = \frac{(x+1)(x^2 - 7)}{x+2}$$

$$f(x) = x^2 - x - 5 + \frac{3}{x+2}.$$

If the degree or the coefficients of $f$ are important, then the first version is appropriate; the second version makes it easy to solve the equation $f(x) = 0$; the last version is a preparation for integrating $f$.

EXAMPLE. Consider the following definition of a subset of the rationals

$$A = \left\{ y \in \mathbb{Q} : y = \frac{x}{x^2 + 1}, \ x \in \mathbb{Z}, \ x < 0 \right\}.$$

The volume of notation is disproportionate to such a simple object. To economise symbols, we switch from the Zermelo to the standard definition of a set, and consider only elements of the required form. We also remove the inequality, and replace it with a negative sign. Finally, we use a more appropriate symbol for the natural numbers.

$$A = \left\{ \frac{-n}{n^2 + 1} : n \in \mathbb{N} \right\}.$$

Now the formula is much more transparent.

EXAMPLE. To simplify the appearance of a complex formula, apply the principle of 'divide and conquer'. In the cluttered formula

$$R = x(ad - bc) - y(ad - cb)^2 + z(ad - cb)^3$$

the expression $ad - bc$ appears as a unit. We exploit this fact to improve the layout

$$R = x\delta - y\delta^2 + z\delta^3 \qquad\qquad \delta = ad - bc.$$

The formula is tidier, and the structure of $R$ is clearer.

EXAMPLE. Our next challenge is to improve an intricate double sum:

$$z(y_1, y_2, \ldots) = \sum_{i=1}^{\infty} \sum_{y=0}^{y_i-1} (y+1)x^{i-1}. \tag{5.10}$$

The meaning of $z$ is not at all evident. This quantity depends on $x$, but its dependence does not appear explicitly; the poor choice of symbols obscures matters further. We note that the parameters $y_i$ are integers, since they are the upper limit of the inner summation; accordingly, we replace the symbol $y$ with $n$ —see (5.1). Then we adopt the 'divide and conquer' principle, splitting up the sum as follows

$$z(x, n_1, n_2, \ldots) = \sum_{i=1}^{\infty} d_i x^{i-1} \qquad d_i = \sum_{k=0}^{n_i-1} (k+1).$$

We now see that $z$ is a power series in $x$; its coefficients are finite sums, determined by the elements of an integer sequence. These are sums of arithmetic progressions, which can be evaluated explicitly

$$\sum_{0 \leqslant k \leqslant n-1} (k+1) = \sum_{1 \leqslant k \leqslant n} k = \frac{n(n+1)}{2}.$$

(In this passage we have dropped the subscript $i$, since the association $n \leftrightarrow n_i$ is clear, and we have switched to the standard sigma-notation to change variable — cf. equation (5.8).) Our original double sum (5.10) can now be written as

$$z(x, \mathbf{n}) = \frac{1}{2} \sum_{i=1}^{\infty} n_i(n_i+1)x^{i-1} \qquad \mathbf{n} = (n_1, n_2, \ldots)$$

where we have used again the 'divide and conquer' method. The dependence of $z$ on the variable $x$ and the sequence $\mathbf{n}$ is now clear.

## 5.5   Writing definitions

A definition requires a pause, to give the reader time to absorb it. This is achieved by giving the definition twice, first with words, then with symbols (or vice-versa), or by using two different formulations, or by supporting the definition with an example. We illustrate some of the possibilities.

EXAMPLE.

*Let $\mathscr{P}$ be a predicate over the integers, that is, a function of the type*

$$\mathscr{P} : \mathbb{Z} \to \{T, F\}.$$

The second part of the sentence reminds the reader of the meaning of predicate over the integers.

EXAMPLE.

*We consider the set A of all co-prime pairs of positive integers, namely*

$$A = \{(m, n) \in \mathbb{N}^2 : \gcd(m, n) = 1\}.$$

The second part of the sentence repeats the definition in symbols, and also reminds the reader of the meaning of co-prime. The Zermelo definitions of a set —cf. (2.7)— should be used sparingly, and should not be used at all if we write for a non-mathematical audience. (Lars Ahlfors in his beautifully written text 'Complex Analysis' deliberately avoids them.)

EXAMPLE.

*Let $\lambda$ be a real number, and let $\Pi(\lambda)$ be the plane in three-dimensional euclidean space, which is orthogonal to the vector $v(\lambda) = (1, \lambda, \lambda^2)$. Thus $\Pi(\lambda)$ consists of all points $z = (x_1, x_2, x_3) \in \mathbb{R}^3$ for which the scalar product $z \cdot v = x_1 + x_2 \lambda + x_3 \lambda^2$ is equal to zero.*

In the first sentence, the dependence of $\Pi$ and $v$ on $\lambda$ is made explicit, which is helpful. The second sentence states the connection between orthogonality and scalar product, and establishes some notation.

EXAMPLE.  Consider the following definition:

*Let $\mathscr{N}$ be the set of sequences of natural numbers, such that every natural number is listed infinitely often. For example, the sequence*

$$(1, 1, 2, 1, 2, 3, 1, 2, 3, 4, \ldots)$$

*belongs to $\mathscr{N}$.*

A non-experienced reader will have no idea that such a construction is at all possible; so we have given an example.

The definition of a symbol should appear as near as possible to where the symbol is first used; defining a symbol immediately after its first appearance is also acceptable, provided that the definition is given *within the same sentence.*

EXAMPLE. We give the same definition three times, articulating the changes in emphasis that accompany each version.

*Consider the power series*

$$h(x) = \sum_{n=1}^{\infty} a_n x^n,$$

*where the coefficient $a_n$ is the square of the $n$-th triangular number.*

The definition of $a_n$ immediately follows its appearance. The displayed formula represents a general power series, and its specific nature is revealed only by reading the entire sentence. For this reason, this format may not be ideal if the formula is to be referred to from elsewhere in the text. This definition puts some burden on the readers who are unfamiliar with the term **triangular number**.

We consider a variant of the definition above.

*Let $t_n$ be the $n$-th triangular number. We consider the power series*

$$h(x) = \sum_{n=1}^{\infty} t_n^2 x^n.$$

Now $t_n$ is defined before being used, with the symbol $t$ chosen so as to remind us of 'triangular'. The formula is clearer. Just glancing at it makes us want to find out what $t_n$ is, and to do this one would begin to scan the text preceding, rather than following, the formula.

Our third version combines verbal and symbolic definitions.

*We consider the power series $h(x)$ whose coefficients are the square of the triangular numbers, namely,*

$$h(x) = \sum_{n=1}^{\infty} t_n^2 x^n \qquad t_n = \frac{n(n+1)}{2}.$$

The formula has become more complex, because all quantities are defined within it. In return, the formula is self-contained; if we number it, we will be able to refer to it from elsewhere in the text.

## 5.6   Introducing a new concept

The sentence that introduces a new concept should appear at the beginning of a paragraph, and the concept should be placed in a prominent position within that sentence. Suppose we want to introduce the logarithm.

BAD:  An important example of a transcendental function is the logarithm.

This is a classic bad opening: '*an example of something is something else*'.  As written, the focus of attention is the transcendental functions, not the logarithm, so we must re-arrange words.

GOOD:  The logarithm is an important example of a transcendental function.

In the next example we want to stress the link between orthogonality of vectors and the vanishing of their scalar product:

BAD:  A commonly used method to check the orthogonality of two non-zero vectors
       is to compute their scalar product and then verify that it is zero.

Besides the dreadful '*a commonly used method for this is that*', this sentence gets lost in describing the computation.

GOOD:  If the scalar product of two non-zero vectors vanishes, then the vectors are
        orthogonal.

GOOD:  If the scalar product of two vectors is zero, then the vectors are orthogonal
        (unless one of the vectors is zero).

Placing a definition at the very beginning of a section may not be appropriate, particularly if we are writing for non-experts. In the following examples, we illustrate some techniques for setting the scene for a definition. We begin by asking a question. This rhetorical device engages the readers, and prepare them for what is about to come.

EXAMPLE.  Introducing recursive sequences. (See section 8.4.)

*Let $n$ be a natural number.  How should we define $2^n$?  We could use repeated multiplication*

$$2^n := \underbrace{2 \times 2 \times \cdots \times 2}_{n},$$

*but we could also write*

$$2^1 := 2 \qquad \text{and} \qquad 2^n := 2 \times 2^{n-1} \quad n > 1.$$

*The second formula is an example of a* **recursive definition** *of a sequence* $(a_n)$. *When* $n = 1$, *the first term* $a_1 = 2$ *is defined explicitly; then, assuming that* $a_{n-1}$ *has been defined, we define* $a_n$ *in terms of it.*

[*The recursive definition of a general sequence follows.*]

The exposition begins with a question, leading to the definition of integer exponentiation in terms of multiplication. The same object is then defined recursively. Due to the simplicity of the example, the exposition is concrete and accessible. After this preamble, the reader will be able to grasp an abstract definition.

EXAMPLE. Introducing the exponential function.

*The process of differentiation turns a real function into another real function. For example, differentiation turns the sine into the cosine. Are there functions that are not changed at all by differentiation?*

[*The definition of the exponential function follows.*]

First, a structural property of differentiation is recalled. The example that follows supports the statement just made, and lead us to a question. The rest of the argument will now develop around the search for functions which are not changed by differentiation.

EXAMPLE. Introducing the rational numbers.

*A rational number is represented by a pair of integers, the numerator and the denominator. Because these integers need not be co-prime, we can choose them in infinitely many ways. Formally, how are we to construct a single rational numbers from an infinite set of pairs of integers? How do we define the set* $\mathbb{Q}$ *from the cartesian product* $\mathbb{Z} \times (\mathbb{Z} \smallsetminus \{0\})$?

[*The definition of the set of rational numbers follows.*]

The first two sentences recall elementary facts. Then a question leads the reader to think about this problem more carefully. A second question, which echoes the first, uses proper terminology and notation, in preparation for a formal construction.

EXAMPLE. We set an exercise, which consists of a list of related questions. The topic is the number of relations on a finite set.

> *What is a relation on a set?*
>
> *Can a relation be defined on the empty set?*
>
> *How many relations can one define on a two-element set?*
>
> *Let n be the number of relations on a set. What values can n assume?*

The first question checks background knowledge; the other questions gently explore the problem, guiding the reader from the specific to the general. Textbook exercises are sometimes structured in this way, to encourage independence in learning. Formulating questions is not merely a device for structuring exercises. Asking the right questions —-those which chart the boundaries of our knowledge— is the essence of research.

## 5.7 Writing a short description

Writing a synopsis of a mathematical topic is a common task. This could be a summary of a section in a book, an abstract, an informal explanation of the content of a theorem. It could be the synthesis of a whole subject, in the introduction or in the closing section of a large document.

Writing a short essay is difficult —the shorter the essay, the greater the difficulty— so this task will be appropriate for the final stages of a learning process. We shall adopt the format of a MICRO ESSAY: 100-150 words (one or two paragraphs) and no mathematical symbols. Within such a confined space, one is forced to make difficult decisions on what to say, and what to leave out; the lack of symbols gives further prominence to the concepts. Our command of the syntax will be put to the test.

Our first MICRO ESSAYis a summary of section 2.1 on sets [⫲, 150]. We have access to all relevant material, and the main difficulty is to decide what are the highlights of that section. We select two ideas: how to define a set, and how to construct new sets from old ones.

*A set is a collection —finite or infinite— of distinct mathematical objects, where the ordering of the objects is immaterial. A small set may be defined by listing explicitly all its elements; a large set is normally defined by specifying the characteristic properties of its elements.*

*With arithmetical operations (sum, multiplication, etc,) we combine numbers to obtain new numbers. Likewise, we may combine sets via set operators (union, intersection, difference, etc), to obtain new sets. Another common construct is the cartesian product of sets, whereby a new set is constructed by assembling together existing sets. A well-known example is the cartesian plane, which is made of two copies of the real line, each representing a coordinate.*

To describe set operators we have used an analogy with arithmetical operators; accordingly, the first two sentences of the second paragraph have the same structure. For the technical term 'cartesian product of sets', we have avoided an exact definition (not enough space!), opting instead for a pictorial description, supported by an illustrative example.

Next we write a MICRO ESSAY on prime numbers, a synthesis of our knowledge of this topic. **[∉, 100]**

*A prime is a positive integer divisible only by itself and unity (however, 1 is not considered prime). The importance of primes in arithmetic stems from the fact that every integer admits a unique decomposition into primes. The infinitude of primes (known from antiquity), and their unpredictability make them object of great mathematical interest.*

*Many properties of an integer follow at once from its prime factorization. For instance, looking at the exponents alone, one can determine the number of divisors, or decide if an integer is a power (i.e.,* [3] *a square or a cube).*

*Primality testing and prime decomposition are computationally difficult problems, that are theoretically challenging and very relevant to applications.*

This essay begins with a definition. The technical point concerning the primality of 1 has been confined within parentheses, to avoid cluttering the first sentence.

---

[3] 'abbreviation for the Latin *id est,* which means 'that is'.

The next two sentences deliver core information in a casual —yet precise— way. We state two important theorems (the Fundamental Theorem of Arithmetic, and Euclid's theorem on the infinitude of the primes) without mentioning the word 'theorem'. We also give a hint of why mathematicians are so fascinated by primes. The second paragraph elaborates on the importance of unique prime factorisation, by mentioning two applications without details. The short closing paragraph echoes the last sentence of the first paragraph, to stimulate the reader's curiosity to learn more.

Now a real challenge: write a MICRO ESSAY on theorem 3.2, page 57, which we reproduce here for convenience.

**Theorem.** *Let $X$ be a set, let $A, B \subseteq X$, and let $\mathscr{P}_A$, $\mathscr{P}_B$ be the corresponding characteristic functions. The following holds (the prime denotes taking complement)*

$$
\begin{array}{rrcl}
(i) & \neg\, \mathscr{P}_A & = & \mathscr{P}_{A'} \\
(ii) & \mathscr{P}_A \wedge \mathscr{P}_B & = & \mathscr{P}_{A \cap B} \\
(iii) & \mathscr{P}_A \vee \mathscr{P}_B & = & \mathscr{P}_{A \cup B} \\
(iv) & \mathscr{P}_A \Rightarrow \mathscr{P}_B & = & \mathscr{P}_{(A \smallsetminus B)'} \\
(v) & \mathscr{P}_A \Leftrightarrow \mathscr{P}_B & = & \mathscr{P}_{(A \cap B) \cup (A \cup B)'} \cdot
\end{array}
$$

This list of inscrutable formulae looks daunting; we must extract a theme from it. We inspect the part of section 3.2 leading to the statement of this theorem: it deals with characteristic functions, and the main idea is to link characteristic functions to sets. We can see such a link in every formula: on the left there are logical operators, on the right set operators. Given the specialised nature of this theorem, we'll have to write for a mathematically mature audience, and a minimum of jargon will be unavoidable.

*The action of a logical operator ( AND, OR, etc.) on a characteristic function generates another characteristic function. Since every characteristic function corresponds to a set, this process is mirrored in new sets being generated from old ones by means of set operations. The theorem establishes a correspondence between these two classes of objects. Under this correspondence, the negation operator, which is boolean, is represented by the set operation of taking the complement. More precisely, the negation of the characteristic function of a set is the characteristic function of the complement of this set. Analogous results are established with respect to the main logical operators.*

Words can't compete with formulae in the delivery of the details. To avoid tedious repetitions, we have chosen to explain just one formula (the easiest one!) carefully, while the other formulae are mentioned under the generic heading 'analogous results'.

## Exercises

**Exercise 5.1.** Consider the following question:

> *Why is it that when the price of petrol goes up by 10% and then comes down 10%, it doesn't finish up where it started?*

1. Write an explanation for the general public. Do not use mathematical symbols, as most people find them difficult to understand. [∉]

2. Write an explanation for mathematicians, combining words and symbols so as to achieve maximum clarity. You should deal with the more general problem of two opposite percentage variations of an arbitrary non-negative quantity (i.e., not specialised to petrol, or to 10% variation, or to the fact that the decrease followed the increase and not the other way around).
   [*Make sure that the variation does not cause the quantity in question to become negative.*]

**Exercise 5.2.** Consider the following question:

> *I drive ten miles at 30 miles an hour, and then another ten miles at 50 miles an hour. It seems to me my average speed over the journey should be 40 miles an hour, but it doesn't work out that way. Why not?*

Write an explanation for the general public, clarifying why such a confusion may arise. You may perform some basic arithmetic, but do not use symbols, as most people find them difficult to understand. [∉]

**Exercise 5.3.** Consider the following question:

> *I tossed a coin four times, and got head four times. It seems to me that if I toss it again I am much more likely to get tail than head, but it doesn't work out that way. Why not?*

Write an explanation for the general public, clarifying why such a confusion may arise. You may use symbols such as $H$ and $T$ for head-tail outcomes, but avoid using other symbols. [∉]

**Exercise 5.4.** Consider the following question:

> *In a game of chance there are three boxes: two are empty, one contains money. I am asked to choose a box, by placing my hand over it; if the money is in that box, I win it. Once I have made my choice, the presenter —who knows where the money is— opens an empty box and then gives me the option to reconsider. I can change box if I wish. It seems to me that changing box would make no difference to my chances of winning, but it does not work out that way. Why not?*

Write an explanation for the general public, explaining what is the best winning strategy. [∉]

**Exercise 5.5.** Formulate a list of questions to help a student approach the following problem:

> *Prove that for all integers $m, n$, we have $m\mathbb{Z} + n\mathbb{Z} = \gcd(m, n)\mathbb{Z}$.*

The assignment should consist of four of five questions of increasing difficulty, the last of which dealing with the statement above.

**Exercise 5.6.** Write a MICRO ESSAY on each of the following topics. [∉, 150]

1. *Quadratic equations and complex numbers.*

2. *Definite vs. indefinite integration.*

3. *Differential and difference equations.*

4. *What are matrices useful for?*

5. *Images and inverse images of sets.*

6. *From pairs of integers to rational numbers.*

7. *Three ways of defining the exponential function.*

**Exercise 5.7.** Write an essay on the following topic.

> *In probability, a random variable is not random and is not a variable: it's a function!*

# Chapter 6

# Forms of argument

A mathematical theory begins with a collection of **axioms** or **postulates**. These are statements that are assumed to be true —no verification or justification is required. The axioms are the building blocks of a theory. From them we deduce other true statements, which are called **theorems.** The process of deduction that establishes a theorem is called a **proof**. As more and more theorems are proved, the theory is enriched by a growing list of true statements.

When developing a mathematical argument, we need to put the sentences in the right order, so that every sentence is either an axiom or a true statement derived from axioms or earlier statements. In practice, only significant statements will be called theorems. In the process of proving a theorem, many true statements may be derived, but these fragments are not usually assigned any formal label. Sometimes a proof rests on statements which, if substantial, are not of independent significance: they are called **lemmas**. Finally, a **proposition** is a statement which deserves attention, but which is not sufficiently general or significant to be called a theorem.

In this chapter we survey some methods to give shape to a mathematical argument. Our survey will continue in chapter 7, which in devoted to inductive arguments.

## 6.1   Anatomy of a proof

A first analysis course begins with a list of axioms defining the real number system. Then one may be asked to prove statements such as

$$\forall x \in \mathbb{R}, \ -(-x) = x. \tag{6.1}$$

These requests are often met with bewilderment: isn't this statement obvious? In recognising the truth of this assertion, we make implicit use of a lot of knowledge, which we derive from our experience, but which is not part of the axioms. A proof from axioms involves erasing all previous knowledge.

We now introduce some axioms, state theorem equivalent to (6.1), and then prove it from the axioms. We put the proof under X-rays, to expose every detail. Our purpose here is to dissect a mathematical argument, not to write an elegant proof. Accordingly, we shall use the language in a mechanical way, articulating every step in the argument. To avoid making implicit assumptions, we represent familiar object with an unfamiliar notation.

*We are given a set $\Omega$, and a binary operator '$\odot$' on $\Omega$, with the following properties*

G1:  $\forall x, y, \in \Omega, \ x \odot y \in \Omega$

G2:  $\forall x, y, z \in \Omega, \ (x \odot y) \odot z = x \odot (y \odot z)$

G3:  $\exists \lozenge \in \Omega, \ \forall x \in \Omega, \ x \odot \lozenge = \lozenge \odot x = x$

G4:  $\forall x \in \Omega, \ \exists x' \in \Omega, \ x \odot x' = \lozenge$

Setting $\Omega = \mathbb{R}$, and '$\odot$' = '$+$', one recognises that $\lozenge$ represents 0 and $x'$ represents $-x$. (These axioms define a very general object: a **group**.) We are ready to state and prove our theorem.

1. **Theorem**.

2. $\forall x \in \Omega, \ x'' = x.$

3. PROOF.

4. Let $x \in \Omega$ be given;

5. then (by G4) $x'' \in \Omega$;

6. let $a := (x \odot x') \odot x'', \ \ b := x \odot (x' \odot x'')$;

7. hence (by G1 and 5) $a, b, \in \Omega$, and (by G2) $a = b$;

8. then (by G3 and G4) $a = \lozenge \odot x'' = x''$;

9. then (by G3 and G4) $b = x \odot \lozenge = x$;

10. hence (by 7, 8, 9) $x'' = x$.

11. □

Some items in this list are about the way we are writing the proof, not about the mathematics itself. Items 1, 3 and 11 locate the statement of the theorem, and the beginning and the end of the proof.

Items 5, 8, 9 begin with the adverb 'then', which we use to say that what comes after is deduced from an axiom. An expressions that says things about the text is called a **logical tag**. Other tags serve the same purpose, for instance 'and therefore'.

Formally, item 5 is our first **theorem**, namely a true statement deduced from the axioms.

Items 7 and 10 are deduction of a slightly different nature, and we have used the tag 'hence' to flag them. These items use facts assembled from previous statements, as well as axioms.

Items 4 and 6 are not statements but **instructions**, in this case to define something. Many different types of instruction can be found in a proof, for instance instructions to draw a picture, or to make an assumption.

In item 10 we use implicitly the **transitivity** of the equality operator (if $x = y$ and $y = z$, then $x = z$). Formally, an axiom is needed to justify this deduction.

Finally, in item 11, the tag '□' signals the end of the proof. An alternative is the acronym Q.E.D, from the Latin *quid erat demonstrandum* which means 'which is what was to be demonstrated'.

In our second example, we look at the proof of the irrationality of $\sqrt{2}$. We use the same style and conventions as in the previous theorem. However, this proof is quite removed from the axioms of arithmetic, in that it uses several definitions and facts without justification. In this respect this proof is more typical than the previous one.

1. **Theorem**.

2. $\sqrt{2}$ is irrational.

3. PROOF.

4. Suppose $\sqrt{2}$ is rational;

5. then $\exists\, m, n \in N$ such that $\sqrt{2} = m/n$.

6. Choose $m$, $n$ in 5 so that they are co-prime.

7. Since $\sqrt{2} = m/n$,

8. then $2 = m^2/n^2$;

9. then $2n^2 = m^2$;

10. then $m^2$ is even;

11. hence (by 10 and theorem A) $m$ is even;

12. then $m = 2h$ for some natural number $h$;

13. hence (by 9 and 12) $2n^2 = 4h^2$;

14. then $n^2 = 2h^2$;

15. then $n^2$ is even;

16. hence (by 15 and theorem A) $n$ is even;

17. hence (by 11 and 16) $m$ and $n$ are not co-prime;

18. hence (by 6 and 17) assumption 4 is false;

19. then its negation is true.

20. $\square$

This proof contains new elements. There are definitions and facts used without justification:

the definition of the rational numbers and of $\sqrt{2}$ (items 2,4,5);

the definition of co-primality (item 6);

some properties of equations (items 8,9);

the definition of an even integer (items 12, 17);

a theorem of arithmetic (items 11,16). We called it 'theorem A', do you know what it is?

The core of the proof is an **assumption** (item 4), where we are instructed to assume that $\sqrt{2}$ is rational. This is puzzling: how can we assume that $\sqrt{2}$ is rational when we know it isn't? How can we let $\sqrt{2} = m/n$ when we haven't been told what numbers $m$ and $n$ are?

This is not so strange; assumptions belong very much to common reasoning.

*Suppose that I have* $100,000,000$ *dollars in a Swiss bank.*

*Suppose that I meet Einstein on top of Mount Everest.*

We are clearly free to explore the logical consequences of these assumptions. In mathematics, assumptions are things we do formally, with symbols on paper rather than in our heads, and according to certain rules. Assumptions are handled using the implication operator $\Rightarrow$ developed in chapter 3.

The assumption that $\sqrt{2}$ is rational, eventually leads to a **contradiction**: items 6 and 17 are conflicting statements. From this fact we deduce that the assumption in item 4 is false (item 18). Here we abandon the assumption 4. But if the statement $\sqrt{2} \in \mathbb{Q}$ is false, then its negation is true. This is item 19, which is what we wanted to prove.

We now abandon the robotic, over-detailed proof style adopted in this section, and analyse various methods of proofs.

## 6.2 Proof by cases

We begin with a simple observation. Sometimes an argument is made tidier by breaking it into a number of cases, precisely one of which holds, and each leading to the desired conclusion.

EXAMPLE. The presence of a piecewise-defined function (see section 4.6) invariably leads to a proof by cases.

**Theorem.** *Let*

$$f : X \to Y \qquad x \mapsto \begin{cases} f_1(x) & \text{if } x \in X_1 \\ f_2(x) & \text{if } x \in X_2. \end{cases}$$

*Prove that, for all $x \in X$, ...*

We have limited information about this statement. From the way the function is defined, we infer that $\{X_1, X_2\}$ is a **partition** of $X$, and that $f_{1,2} : X_{1,2} \to Y$. The statement to be proved contains a universal quantifier.

PROOF. Let $x \in X$ be given. We have two cases

*Case I: $x \in X_1$. Then $f = f_1$, ...*
*Case II: $x \in X_2$. Then $f = f_2$, ...*

The proof begins with a standard opening sentence ('Let $x \in X$ be given'), which acknowledges that $x$ is an arbitrary element of $X$, because it's controlled by a universal quantifier. (Cf. the proof of the first theorem in section 6.1.) With reference to the discussion following expression (3.20), the choice of $x$ is 'my opponent's move'. The second sentence announces that the proof will develop into two cases, which are determined by the piecewise definition of the function $f$. We have shaped the beginning of the proof without knowing what we are supposed to prove!

EXAMPLE. The absolute value function is piecewise defined (see equation (4.12)); its presence leads to proof by cases.

> Prove that the solution set of the inequality $2|x-1| \geqslant |x-2|$ is equal to $(-\infty, 0] \cup [4/3, \infty)$.

PROOF. Let $x$ be a real number. There are three cases.
Case I: If $x < 1$, then the inequality is $2(1-x) \geqslant 2-x$, giving $x \leqslant 0$.
Case II: If $1 \leqslant x < 2$, then the inequality is $2(x-1) \geqslant 2-x$, giving $x \geqslant 4/3$.
Case III: If $x \geqslant 2$, then the inequality is $2(x-1) \geqslant x-2$, which is always verified, giving $x \geqslant 2$.
So the required solution set is the set of $x$ such that $x \leqslant 0$ or $4/3 \leqslant x < 2$, or $x \geqslant 2$, which is the union of two rays

$$(-\infty, 0] \cup [4/3, 2) \cup [2, \infty) = (-\infty, 0] \cup [4/3, \infty),$$

as desired.   $\square$

  Note the careful distinction between strict and non-strict inequalities in the three cases, to avoid missed or repeated values of $x$.

EXAMPLE. A proof by cases, about divisibility.

> Let $n$ be an integer; then $n^5 - n$ is divisible by 30.

PROOF. We express $n^5 - n$ in factored form

$$n^5 - n = n(n^4 - 1) = n(n-1)(n+1)(n^2 + 1).$$

We have the prime factorisation $30 = 2 \times 3 \times 5$. For each prime $p = 2, 3, 5$ we will show that, for any integer $n$, at least one factor of $n^5 - n$ is divisible by $p$.

1.  $p = 2$. Let $n = 2k + j$, for some $k \in \mathbb{Z}$ and $j = 0, 1$.

    If $j = 0$, then $n$ is divisible by 2.

    If $j = 1$, then $n - 1$ is divisible by 2.

2. $p = 3$. Let $n = 3k + j$, for some $k \in \mathbb{Z}$ and $j = 0, \pm 1$.

   If $j = 0$, then $n$ is divisible by 3;

   if $j = \pm 1$, then $n \mp 1$ is divisible by 3.

3. $p = 5$. Let $n = 5k + j$, for some $k \in \mathbb{Z}$ and $j = 0, \pm 1, \pm 2$.

   If $j = 0$, then $n$ is divisible by 5.

   If $j = \pm 1$, then $n \mp 1$ is divisible by 5.

   If $j = \pm 2$, then $n^2 + 1 = 25k^2 \mp 20k + 5 = 5(5k^2 \mp 4k + 1)$ is divisible by 5.

The proof is complete. $\square$

The proof begins by factoring the polynomial $n^5 - n$ and the integer 30. Then we declare our intentions. The expression 'for any integer $n$' acknowledges the presence of the (hidden) universal quantifier $\forall n \in \mathbb{Z}$.

Each prime $p$ leads to $p$ cases, corresponding to the possible values of the remainder $j$ of division by $p$. To simplify the book-keeping, we consider also negative remainders, and then pair together the remainders which differ by a sign. Note the presence of the symbols $\pm$ and $\mp$ within the same sentence. The positive sign in the first expression matches the negative sign in the second expression, and vice-versa.

## 6.3 Implications

Many statements take the form of implications:

> *If P then Q.*

For example:

> *If $p$ is an odd prime, then $2^{p-1} - 1$ is divisible by $p$.*
> *If the sequence $(a_k)$ is periodic, then $(a_{2k})$ is also periodic.*

Implications often appear in disguise. The statements

> *A is a subset of B.*
> *Every repeating decimal is rational.*
> *The determinant of an invertible matrix is non-zero.*

do not contain the implication operator '$\Rightarrow$' or the 'if...then' construct, yet they are implications. This is seen by re-writing them as follows (cf. equation (3.19))

*If $x \in A$, then $x \in B$.*
*If $r$ is a repeating decimal, then $r$ is rational.*
*If $A$ is an invertible matrix, then $\det(A) \neq 0$.*

This formulation makes the implication clear, and also gives us a name to use for the relevant quantities.

Many implications contain a hidden universal quantifier

$$\forall x \in X, \ \mathscr{P}(x) \Rightarrow \mathscr{Q}(x) \tag{6.2}$$

where $\mathscr{P}$ and $\mathscr{Q}$ are predicates over $X$. We rewrite a statement considered above so as to make the quantifier visible

*For all real numbers $r$, if the decimal digits of $r$ are repeating, then $r$ is rational.*

Spelling out an implication in this way is helpful.

### 6.3.1   Direct proof

Let us consider the truth table of the implication operator, given in (3.8). If $P$ is false, then the expression $P \Rightarrow Q$ is true regardless of the value of $Q$, and there is nothing to prove. If $P$ is true, then the implication is true provided that $Q$ is true. So a **direct proof** of the implication amount to **assuming** $P$ and **deducing** $Q$. The assumption of $P$ lasts only until we have reached $Q$. When we have finished proving $Q$, the assumption $P$ is *discharged*, meaning that it can no longer be used.

Every direct proof of *'If P then Q'* must contain a section during which $P$ is assumed. The start of the block is announced by a sentence such as

*Assume P*          *Suppose P*          *Let P.*

The task of the rest of the proof is to prove not the original theorem, but $Q$. We make this clear by writing

*RTP: Q*

where RTP stands for *Remains To Prove* (or *Required To Prove*). The block ends when the proof of $Q$ is completed. This is announced with a closing sentence such as

*We have proved Q*          *The proof of Q is complete.*

The above considerations suggest what should be the first step in a direct proof of an implication

> **Theorem.** *If $p > 3$ is a prime and $p + 2$ is also prime, then $p + 4$ is composite.*
> PROOF. Suppose $p$ is a prime number greater than 3, such that $p + 2$ is prime. RTP: $p + 4$ is composite.

All three assumptions ($p > 3$, $p$ prime, $p + 2$ prime) must be used in the proof, or the proof will be wrong.

Often you can work out how the proof of an implication must start, *even if you haven't the faintest idea of what the maths is about*. We illustrate this point with some real life examples:

> **Theorem 1.** *A closed subset of a compact set is compact.*
> PROOF. Let $X$ be a compact set, and let $C$ be a subset of $X$. Assume that $C$ is closed. RTP: $C$ is compact.

> **Theorem 2.** *If $\lambda \in \mathbb{C}$ is a root of a monic polynomial whose coefficients are algebraic integers, then $\lambda$ is an algebraic integer.*
> PROOF. Let $p$ be a monic polynomial whose coefficients are algebraic integers, and let $\lambda \in \mathbb{C}$ be a root of $p$. RTP: $\lambda$ is an algebraic integer.

> **Theorem 3.** *Every finite basis of a finite-dimensional linear space has the same number of elements.*
> PROOF. Let $V$ be a finite-dimensional linear space, and let $B_1$ and $B_2$ be two finite bases for $V$. Suppose $B_1$ consists of $n_1$ elements and $B_2$ consists of $n_2$ elements. RTP: $n_1 = n_2$.

Establishing a good notation is often decisive, and in all examples above the proof begins by giving names to things. Some authors make this unnecessary by including names in the theorem; others obscure the statement of a theorem by putting too many names in it. We now rewrite the last theorem in such a way as to establish some notation within the statement.

> **Theorem 3.** *Let $V$ be a finite-dimensional linear space. Then every finite basis for $V$ has the same number of elements.*

> **Theorem 3.** *Let $V$ be a finite-dimensional linear space, and let $B_1$ and $B_2$ be two finite bases for $V$. Then $\#B_1 = \#B_2$.*

Let us compare the three formulations of this theorem: the first statement is concise and forceful; the second contains a minimum of notation, which does no harm, but is also unnecessary. The last version establishes some useful notation. In general, this is best done within the proof. However, introducing notation in a theorem is appropriate if one plans to use this notation at a later stage. For instance, one could find the sentence

>  Let $V$, $B_1$, $B_2$ be as in theorem 3

in the material following theorem 3.

## 6.3.2  Proof by contrapositive

In section 3.1 we have seen that the expressions

$$P \Rightarrow Q \qquad \text{and} \qquad \neg Q \Rightarrow \neg P$$

are **equivalent,** namely they are both true or both false for any choice of $P$ and $Q$. The second implication is called the **contrapositive** of the first. It then follows that the first is the contrapositive of the second.

The equivalence between an expression and its contrapositive gives us a method of proving implications, called a **proof by contrapositive,** which is a useful alternative to a direct proof, as sometimes it is easier. To give a proof by contrapositive of

>  If $P$ then $Q$,

we prove instead that

>  If not $Q$ then not $P$,

that is, we start by assuming that $Q$ is false, and then we deduce that $P$ is false.

A proof by contrapositive is structurally identical to a direct proof, only the predicates are different. Thus in place of (6.2) we prove

$$\forall x \in X, \ \neg \mathscr{Q}(x) \Rightarrow \neg \mathscr{P}(x). \tag{6.3}$$

Note that the contrapositive of a quantified implication is obtained by forming the contrapositive of the implication, *without altering the quantifier.* We just replace a boolean function an equivalent function, much like replacing $\sin(x)^2$ with $1 - \cos(x)^2$.

If we have to prove an implication, how do we decide between a direct proof and a proof by contrapositive? There is no absolute criterion; we must compare the assumptions $P$ and $\neg Q$, and decide which of the two is easier to handle. Sometimes it's necessary to try both approaches to find out! We illustrate this point with examples.

EXAMPLE. Consider the statement

$$\forall n \in \mathbb{N}, \ (2^n < n!) \Rightarrow (n > 3) \tag{6.4}$$

The assumption $\mathscr{P}(n) = (2^n < n!)$ is problematic, because its value is not easily computable. By contrast, $\neg\mathscr{Q}(n) = (n \leqslant 3)$ is straightforward, and the contrapositive implication

$$\forall n \in \mathbb{N}, \ n \leqslant 3 \Rightarrow (2^n \geqslant n!) \tag{6.5}$$

involves checking that the boolean expression $(2^n \geqslant n!)$ is true for only three values of $n$:

PROOF. We only have to check three cases:

$$\begin{aligned}
n = 1: & \quad 2 = 2^1 \geqslant 1! = 1 \\
n = 2: & \quad 4 = 2^2 \geqslant 2! = 2 \\
n = 3: & \quad 8 = 2^3 \geqslant 3! = 6.
\end{aligned}$$

Thus the expression (6.5) is true, and the proof is complete.  □

EXAMPLE. Consider the statement

*If the average of four distinct integers is equal to 10, then one of the integers is greater than 11.*

The direct implication involves an assumption on an average value, which entails loss of information; the contrapositive implication involves four integers of bounded size. We opt for the latter, which seems easier:

*Given four distinct integers not greater than 11, their average is not equal to 10.*

PROOF. We prove the contrapositive of the above statement. Let four distinct integers be given. If none of them exceeds 11, then the largest value their sum can assume is $11 + 10 + 9 + 8 = 38$. So the largest possible average is

$$\frac{38}{4} = \frac{19}{2} < 10$$

as desired.  □

## 6.4   Proving conjunctions

A **conjunction** is a statement of the type

>   *P and Q.*

The statements $P$ and $Q$ are called its **conjuncts**. Mathematical statements often take the form of conjunctions; like implications, sometimes conjunctions are hidden.

   A direct proof of the conjunction *P and Q* consists of a proof of $P$ and a proof of $Q$. The two proofs should be separated clearly: it's common practice to put the separate conjuncts in an ordered list, say $i)$ and $ii)$. Then the proof itself should use the same labels $i)$ and $ii)$. It is advisable to begin the proof of each conjunct by announcing our intentions.

PROOF.

   $i)$   We prove $P$. ...

   $ii)$   We prove $Q$. ...

There could be more than two conjuncts; they may be listed using lower-case roman numerals: $i)$, $ii)$, $iii)$, $iv)$, etc,

   Equality of two sets is a hidden conjunction.

>   *The sets A and B are equal.*

What are the conjuncts? Two sets $A$ and $B$ are defined to be equal if they have the same elements. Equivalently, every element of $A$ is an element of $B$, and vice-versa, and so we write

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

Recalling the definition of subset, the structure of a proof is clear.

PROOF.

   $i)$   We prove that $A \subset B$. Let $x \in A$ be given. RTP: $x \in B$.

   $ii)$   We prove that $B \subset A$. Let $x \in B$ be given. RTP: $x \in A$.

   A **bound** is a hidden conjunction.

$$\forall x \in \mathbb{R}, \ |f(x)| < g(x)$$

The conjuncts are

$$\forall x \in \mathbb{R}, \ f(x) < g(x) \qquad \text{and} \qquad \forall x \in \mathbb{R}, \ f(x) > -g(x).$$

The equivalence of two statements is the quintessential conjunction, since the equivalence operator $\Leftrightarrow$ is the conjunct of an implication and its converse —see definition (3.11.

> *P if and only if Q.*

This statement may be rephrased as '*P is necessary and sufficient for Q*', or as $P \Leftrightarrow Q$. The structure of the proof is predictable.

PROOF.

  *i*) We prove $P \Rightarrow Q$. ...

  *ii*) We prove $Q \Rightarrow P$. ...

We could replace either part by a proof of the contrapositive implication.

The following theorem provides an alternative characterisation of primality.

**Theorem** (Wilson). *A natural number $n > 1$ is prime if and only if $n$ divides $(n - 1)! + 1$.*

The theorem says $P \wedge Q$, where

$$P = \forall n > 1, \ (n \text{ is prime}) \ \Rightarrow \ (n \mid (n-1)! + 1)$$
$$Q = \forall n > 1, \ (n \mid (n-1)! + 1) \ \Rightarrow \ (n \text{ is prime}).$$

The outline of the proof is now clear.

PROOF.

  *i*) Let $n$ be a prime number.
      RTP: $n$ divides $(n - 1)! + 1$.

  *ii*) Let $n$ be a natural number greater than 1 which divides $(n - 1)! + 1$.
      RTP: $n$ is prime.

The term 'precisely' may be used to turn a one-sided implication into an equivalence, hence a conjunction.

> *The sine function vanishes at the integer multiples of $2\pi$.*
> *The sine function vanishes precisely at the integer multiples of $\pi$.*

The first statement is a one-sided implication: the set $2\pi\mathbb{Z}$ is contained in the solution set of the equation $\sin(x) = 0$. The second statement is en equivalence: the set $\pi\mathbb{Z}$ is the solution set of that equation.

We restate Wilson's theorem using this expression.

> *The prime numbers are precisely the integers $n$ greater than 1 which divide* $(n-1)!+1$.

## 6.5   Circular arguments

The term 'circular argument' is sometimes used in a negative sense, to characterise a type of mistakes that originate from assuming what we are supposed to prove (section 9.3). Here we use it to characterise chains of implications where the last proposition in the chain coincides with the first one.

For instance, proving the equivalence of two statement $P_1$ and $P_2$, may be thought of as establishing the chain of two implications

$$P_1 \Rightarrow P_2 \Rightarrow P_1.$$

More generally, one considers chains of $n$ implications

$$P_1 \Rightarrow P_2 \Rightarrow \cdots \Rightarrow P_{n-1} \Rightarrow P_n \Rightarrow P_1.$$

Proving all implications in the chain amounts to proving that all statements are equivalent, namely that $P_i \Leftrightarrow P_j$, for all $i, j = 1 \ldots, n$.

For example, let $G$ be a group and $H$ a subgroup of $G$. A **left coset** of $H$ in $G$ is a set $gH = \{gh : h \in H\}$ where $g$ is an element of $G$. A standard result of group theory states that if $x, y$ are elements of $G$ then the following are equivalent:

    (a)   $xH = yH$
    (b)   $xH \subseteq yH$
    (c)   $xH \cap yH \neq \emptyset$
    (d)   $y^{-1}x \in H$.

These are usually proved in a circle, for example $(a) \Rightarrow (b) \Rightarrow (c) \Rightarrow (d) \Rightarrow (a)$, but other arrangements are possible, e.g., $(a) \Leftrightarrow (b) \Rightarrow (c) \Rightarrow (d) \Rightarrow (b)$.

In chapter 7 we will show the equivalence of four formulations of the principle of induction, using a circular argument.

## 6.6  Proof by contradiction

A proof by contradiction of a proposition $\mathscr{L}$ consists of assuming $\neg\mathscr{L}$, and deducing a false statement, i.e., proving the implication

$$\neg\mathscr{L} \Rightarrow \text{FALSE}. \tag{6.6}$$

The false statement to be deduced can be anything, including any assertion that contradicts the assumption $\neg\mathscr{L}$.

To see why contradiction works, suppose we have shown that (6.6) holds. This gives us the boolean equation

$$(\neg\mathscr{L} \Rightarrow \text{FALSE}) = \text{TRUE}$$

to be solved for $\mathscr{L}$. Using the truth table (3.8) of the operator $\Rightarrow$, this equation becomes $\neg\mathscr{L} = \text{FALSE}$, hence $\neg(\neg\mathscr{L}) = \mathscr{L} = \text{TRUE}$, by virtue of (3.7).

EXAMPLE. We consider the formal aspects of the proof by contradiction of the irrationality of $\sqrt{2}$, which we developed in section 6.1. We define three boolean expressions. The first one is

$$P := (\sqrt{2} \in \mathbb{Q}).$$

The value of $P$ is to be determined later. The second expression is

$$Q := \text{`}\exists m, n \in \mathbb{N}, \ (m, n \text{ are co-prime}) \Rightarrow (m, n \text{ are even})\text{'}$$

and finally

$$\mathscr{L} := (P \Rightarrow Q).$$

The expression $Q$ is clearly false, while the steps 4–10 in the proof of the theorem show that $\mathscr{L}$ is true. Then $P$ is false and $\neg P = (\sqrt{2} \notin \mathbb{Q})$ is true, which is what we wanted to prove.

EXAMPLE. A classic proof by contradiction is Euclid's proof of the infinitude of primes.

**Theorem.** *The number of primes is infinite.*

PROOF. Assume there are only finitely many primes, $p_1, \ldots, p_n$. Consider the integer

$$N = 1 + \prod_{k=1}^{n} p_k.$$

Then $N$ is greater than all the primes; moreover, if we divide $N$ by $p_k$ we get remainder 1, and therefore $N$ is not divisible by any of the primes. Now, every integer greater than 1 is divisible by some prime. It follows that $N$ must be divisible by a prime not in the list above, a contradiction.    □

The statement 'every integer greater than 1 is divisible by some prime', which is essential to the proof, needs some justification (see section 7.1).

The proof by contradiction of the implication $P \Rightarrow Q$ takes the form

$$\neg(P \Rightarrow Q) \Rightarrow \text{FALSE}.$$

Now, $\neg(P \Rightarrow Q)$ is equivalent to $P \wedge \neg Q$ (from theorem 3.1.2, page 55) so the above becomes

$$(P \wedge \neg Q) \Rightarrow \text{FALSE}.$$

This form of proof by contradiction is called the **both ends method**: To prove

   If $P$ then $Q$

we assume both $P$ and not-$Q$, and we deduce something impossible. The appealing feature of this method is that it gives us two assumptions, $P$ and $\neg Q$, to exploit. This is often how researchers work, but it easily causes confusion in writing.

We show a poor example of both ends proof, from a textbook:

**Theorem.** *For all real numbers $x$, if $x > 0$, then $1/x > 0$.*

BAD PROOF: If $1/x < 0$ then $(-1)/x > 0$, so $x((-1)/x) > 0$, i.e., $-1 > 0$, a contradiction.  Therefore $1/x \geqslant 0$.  Since $1/x$ can't be 0, we conclude that $1/x > 0$.
□

In this proof, the assumption of $P$ is hidden, and the assumption of $\neg Q$ is confused. Not helpful writing!

PROOF. We prove it by contradiction. Let us assume that $x > 0$ and $1/x \leqslant 0$. Multiplying the second inequality by $-1$, we obtain $-1/x \geqslant 0$, and since $x$ is positive, multiplication by $x$ yields $x(-1/x) \geqslant 0$. Simplification gives $-1 \geqslant 0$, which is false.
□

## 6.7    Counterexamples and conjectures

To prove that a proposition regarding all elements of a set is false, it is sufficient to exhibit a single element of that set for which the proposition fails. This construct is called a **counterexample.**

Formally, to show that the statement

$$\forall x \in X, \ \mathscr{P}(x) \tag{6.7}$$

is false, we prove its negation, namely

$$\exists x \in X, \ \neg\mathscr{P}(x).$$

The existential quantifier implies that it's up to us to find a value of $x$ for which the above expression is true. A single value of $x$ will suffice.

To illustrate the idea of a counterexample, let us consider the integer polynomial $p(n) = n^2 + n + 41$, proposed by Euler. We evaluate $p(n)$ at some integer values of $n$.

| $n$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | $\cdots$ | 20 | $\cdots$ | 30 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $p(n)$ | 41 | 43 | 47 | 53 | 61 | 71 | 83 | 97 | 131 | 151 | 181 | $\cdots$ | 461 | $\cdots$ | 971 |

All these values of $p(n)$ are prime! Moreover, $p(-n) = p(n+1)$, so it seems plausible —if a bit daring— to put forward the following conjecture

**Conjecture 1.** *For all integers n, $p(n)$ is prime.*

A conjecture is a statement that we wish it were a theorem. Three things may happen to a conjecture: *i*) someone produces a proof, and the conjecture becomes a theorem; *ii*) someone produces a counterexample, and the conjecture is proved false; *iii*) none of the above, and the conjecture remains a conjecture.

Our case is *ii*). The negation of conjecture 2 is

    *There is an integer n such that $p(n)$ is composite.*

To prove it, we must exhibit such a value of $n$. For $n = 40$, we find

$$p(40) = 40^2 + 40 + 41 = 40(40+1) + 41 = 41(40+1) = 41^2.$$

So $p(40)$ is not prime, and this counterexample shows that conjecture 1 is false.

For a mathematician, the fear of a counterexample accompanies the formulation of a conjecture. Let us return to conjecture 1: we know it's false, but can it be modified into a weaker statement? It can be verified that $n = 40$ is the smallest value for $n$ for which $p(n)$ is not prime[1]. Are there other such values? If these values were exceptional, a meaningful weaker version of conjecture 1 could be

**Conjecture 2.** *For all but finitely many integers n, $p(n)$ is prime.*

Even if $p(n)$ were composite for, say, a billion values of $n$, the conjecture would still hold. Unfortunately, conjecture 2 is not true either. Its negation

---

[1]This phenomenon has deep roots, see [4, page 155].

*There are infinitely many integers n for which p(n) is composite*

is established by the simple, devastating counterexample

$$p(41k) = (41k)^2 + 41k + 41 = 41(41k^2 + k + 1) \qquad k = 1, 2, \ldots . \qquad (6.8)$$

which says that 41 is a proper divisor of $p(41k)$ for infinitely many values of $k$.

Let's make one final attempt to savage something from our original claim.

**Conjecture 3.** *There are infinitely many integers n such that p(n) is prime.*

This is a very meaningful conjecture. Nobody has ever been able to prove that any quadratic polynomial with integer coefficients assumes infinitely many prime values. On the other hand, if we perform a numerical experiment we discover that, as $n$ increases, $p(n)$ continues to provide a large supply of prime values, even though the composite values slowly become dominant. So this conjecture is very likely to be true, and it makes sense to put it forward.

Arguably, the most famous conjecture in mathematics is the **Riemann's hypothesis** (RH), formulated in 1859. This conjecture (essentially) says that the zeta-function defined in equation (5.3), page 101 vanishes only at points whose imaginary part is equal to $1/2$. The depth of this statement is not at all apparent from its seemingly inconspicuous formulation.

In a sense, the mathematics community cannot afford to wait for this important matter to be settled, and several theorems about the RH have been proved. Some theorem formulate the RH in terms of equivalent statements, establishing the importance of the RH in many areas of mathematics. There are also proofs of **conditional theorems**, which assume the validity of the RH, and deduce some of its consequences.

Other famous conjectures include the **twin-prime conjecture,** attributed to Euclid

**Conjecture.** *There are infinitely many primes p such that p + 2 is also prime.*

and the **Goldbach conjecture**

**Conjecture.** *Every even integer greater than 2 can be written as a sum of two primes.*

Both conjectures are easy to formulate, and widely believed to be true, but a proof seems hopelessly difficult. This situation is common in the theory of numbers, where innocuous-looking questions may range —quite unexpectedly— from the very easy to the very difficult. To get a taste of this phenomenon, let us arrange the natural numbers into four columns, with the primes highlighted in boldface:

| 1 | **2** | **3** | 4 |
|---|---|---|---|
| **5** | 6 | **7** | 8 |
| 9 | 10 | **11** | 12 |
| **13** | 14 | 15 | 16 |
| **17** | 18 | **19** | 20 |
| 21 | 22 | **23** | 24 |
| 25 | 26 | 27 | 28 |
| **29** | 39 | **31** | 32 |
| 33 | 34 | 35 | 36 |
| **37** | 38 | 39 | 40 |
| **41** | 42 | **43** | 44 |
| 45 | 46 | **47** | 48 |
| 49 | 50 | 51 | 52 |
| **53** | 54 | 55 | 56 |
| 57 | 58 | **59** | 60 |
| **61** | 62 | 63 | 64 |
| 65 | 66 | **67** | 68 |
| 69 | 70 | **71** | 72 |
| **73** | 74 | 75 | 76 |
| ⋮ | ⋮ | ⋮ | ⋮ |

We ask some questions about this table, listed in order of increasing difficulty.

*1. Why is there just one prime in the second column, and no primes at all in the fourth?*

*2. The first two rows contain four primes in total; are there other pairs of adjacent rows with the same property?*

*3. How many primes does the table contain?*

*4. Are there infinitely many rows with no primes in them?*

*5. Are there adjacent rows with no primes? What about sequences of three or more adjacent rows with no primes?*

*6. Are there infinitely many primes in the first column?*

*7. Are 'half' of the primes in column 1 and 'half' in column 3?*

*8. Are there infinitely many rows containing two primes?*

At the end of primary school, a pupil should be able to answer question 1, while in secondary school one could see why the answer to question 2 must be negative

(one of three consecutive odd integers is divisible by 3).  The answer to questions
3 to 7 is affirmative.  Euclid's theorem (page 133), which answers question 3, is
normally taught in a first year university course, although we have seen that the
proof does not require advanced ideas.  Answering questions 4 and 5 will involve
some ingenuity, but no major difficulty (see exercises).

The mathematics needed to answer 6 and 7 becomes difficult.  The proof that
column 1 contains infinitely many primes, known since the 17th Century, is given
today in an undergraduate course in number theory.  The formulation of question 7
can be made precise, and the affirmative answer was given at the end of the 19th
Century.  Understanding the proof requires a solid background in analysis.

Many mathematicians would conjecture that the answer to question 8 is again
positive, but at present nobody knows how to prove it.  This is a variant of the twin
primes conjecture stated above.

## 6.8   Writing a good proof

Proofs come in all shapes and sizes.  Correctness is, of course, imperative, but a good
proof requires a lot more.  One must realise that reading a proof is demanding, and
that any assistance will be welcome.  At the same time, there may be the conflicting
requirement to keep the proof's length within manageable bounds.  The author's
style, and the mathematical maturity of the target audience will also be significant.

There is one universal rule.  At key junctures in a proof:

- STATE CONCISELY WHAT YOU PLAN TO DO.

- WHEN YOU'VE DONE IT, SAY SO.

We now examine statements of theorems and proofs that are less than optimal,
and turn them into clear statements and good proofs.  Our proofs are very detailed,
and suitable for beginning university students.  Occasionally, we provide a concise
version of a proof, addressed to a mathematically mature audience.

EXAMPLE.  The following theorem is taken from a textbook

BAD THEOREM.  *If $x^2 \neq 0$, then $x^2 > 0$.*

BAD PROOF.  If $x > 0$ then $x^2 = xx > 0$.  If $x < 0$ then $-x > 0$ , so $(-x)(-x) > 0$, i.e.,
$x^2 > 0$.   □

The statement of the theorem is incomplete, in that there is no information about
the quantity $x$.  The proof is rather concise.  The theorem is an implication, but what

has happened to the assumption $x^2 \neq 0$? If the assumption is obvious, one may leave it out, but one must judge whether the readers are sophisticated enough to see what it must be. Furthermore, this is a proof by cases (see section 6.2), and it would be helpful to make this more explicit.

THEOREM. *For all real numbers $x$, if $x^2 \neq 0$, then $x^2 > 0$.*

PROOF. Let $x$ be a real number such that $x^2 \neq 0$. Then $x \neq 0$, and have two cases:
$i)$ $x < 0$. Then $-x > 0$, so $(-x)(-x) > 0$, that is, $x^2 > 0$.
$ii)$ $x > 0$. Then $x^2 = xx > 0$. $\quad\square$

EXAMPLE. We consider an arithmetical statement.

BAD THEOREM. $\forall n \in \mathbb{Z}$, $2 \nmid n \Rightarrow 8 | n^2 - 1$.

BAD PROOF.
$2 \nmid n \Rightarrow \exists k \in \mathbb{Z} \; n = 2k + 1$;
$n^2 - 1 = 4k(k+1)$;
$\forall k \in \mathbb{Z}$, $2 | k(k+1)$. $\quad\square$

Both statement and proof are unnecessarily formal. This theorem is best stated with a mixture of words and symbols.

THEOREM. *If $n$ is an odd integer, then $n^2 - 1$ is divisible by 8.*

PROOF. Let $n$ be an odd integer. We will show that $n^2 - 1 = 8j$, for some $j \in \mathbb{Z}$.
Since $n$ is odd, we have $n = 2k + 1$, for some integer $k$. We find

$$n^2 - 1 = (2k+1)^2 - 1 = 4k^2 + 4k + 1 - 1 = 4k(k+1) \tag{1}$$

which shows that $n^2 - 1$ is divisible by 4. Now, one of $k$ or $k + 1$ is even, and therefore their product is even. Thus we have $k(k+1) = 2j$, for some $j$. Inserting this expression in (1), we find $n^2 - 1 = 8j$, as desired. $\quad\square$

The proof's opening sentence acknowledges the presence of a hidden universal quantifier ($\forall n \in \mathbb{Z}$). What needs to be done is then clarified. Checking divisibility by 2 is a proof by cases, but there's no need to announce it, given its simplicity.

The given proof is quite detailed, and it can safely be shortened.

CONCISE PROOF. An odd integer $n$ is of the form $n = 2k + 1$, for some integer $k$. A straightforward manipulation gives $n^2 - 1 = 4k(k+1)$. Our claim now follows from the observation that the product $k(k+1)$ is necessarily even. $\quad\square$

The adverb 'necessarily' is inserted to signal that the conclusion (the product $k(k+1)$ is even) requires a moment's thought. The expression 'straightforward manipulation' is used to omit some steps in the derivation, while warning the reader that

some calculations are required. In such a circumstance, the adjective 'straightfor-
ward' is preferable to 'easy', which is subjective, or 'trivial', which carries a hint of
arrogance. For instance, it's quite appropriate to claim that the arithmetical identity

$$29 \;=\; \left(2+\frac{\sqrt{-1}\,(1+\sqrt{5})}{2}\right)\left(2-\frac{\sqrt{-1}\,(1+\sqrt{5})}{2}\right)$$
$$\times \left(2+\frac{\sqrt{-1}\,(1-\sqrt{5})}{2}\right)\left(2-\frac{\sqrt{-1}\,(1-\sqrt{5})}{2}\right)$$

can be verified with a 'straightforward calculation'.

EXAMPLE. A relation between rational and irrational numbers.

BAD THEOREM. $\forall a,b \in \mathbb{R}, \; (a \in \mathbb{Q} \wedge b \notin \mathbb{Q}) \Rightarrow a+b \notin \mathbb{Q}$.

BAD PROOF: Spse $a+b \in \mathbb{Q}$. Then $a+b = m/n$.
If $a \in \mathbb{Q}$, then $a = p/q$, and $b = m/n - p/q \in \mathbb{Q}$.   $\square$

The symbolic formulation obscures the simple content of the theorem. The
proof is also simple, but it introduces unnecessary symbols, while clarity could be
improved with some comments.

THEOREM. *The sum of a rational number and an irrational number is irrational.*

PROOF. Let $a$ and $b$ be a rational and an irrational number, respectively. Consider
the identity $b = (b+a) - a$. If $a+b$ were rational, then $b$ would also be rational,
being the difference of two rational numbers. This contradiction shows that $a+b$
must be irrational.   $\square$

EXAMPLE. In what follows, everything (definition, statement of theorem, and proof)
need re-writing.

BAD DEFINITION. Let $b_j$ be decimal digits. We define

$$\overline{b_1 \cdots b_n} = b_1 \cdots b_n b_1 \cdots$$

BAD THEOREM. *Let $m \in \mathbb{Z}$, and let $x = m + 0.a_1 \cdots a_k \overline{b_1 \cdots b_n}$. Then $x \in \mathbb{Q}$.*

BAD PROOF: Let $A = 0.a_1 \cdots a_k$, $B = 0.b_1 \cdots b_n$. Then $A, B \in \mathbb{Q}$, and

$$0.\overline{b_1 b_2 \ldots b_n} = B\left(1 + \frac{1}{10^n} + \frac{1}{10^{2n}} + \cdots\right) = B\frac{10^n}{10^n - 1} =: B' \in \mathbb{Q}.$$

Thus $x = m + A + B'10^{-k} \in \mathbb{Q}$. $\quad\square$

In the definition, the periodicity of digit sequence is not obvious. In the statement of the theorem, the digits $a_j$ are undefined. In addition, a hybrid notation is used for $x$, which is the sum of an integer expressed symbolically, and a decimal number.

As we did before, we state the theorem with words, which is more effective and does not require new notation. The main notation is introduced *after* the statement of the theorem, but *before* the beginning of the proof. This arrangement should be considered if the notation is needed outside the confines of the proof, or if the proof is heavy and needs lightening up, or simply as a variation from the rigid definition-theorem-proof format. For the purpose of reference, we number the various steps in the argument.

THEOREM. *Every number whose decimal digits eventually repeat is rational.*

1. Before proving the theorem, we establish some notation. Let $a_1 a_2 \cdots$ and $b_1 b_2 \cdots$ be strings of decimal digits: $a_i, b_i \in \{0, 1, \ldots, 9\}$. We use the over-bar to denote a string of digits consisting of a pattern repeating indefinitely

$$\overline{b_1 \cdots b_n} = \underbrace{b_1 \cdots b_n}_{n} \underbrace{b_1 \cdots b_n}_{n} \cdots. \tag{6.9}$$

PROOF.
2. Let $x$ be a real number with eventually repeating decimal digits. Then $x$ has a decimal representation of the type

$$x = a_0 \cdots a_j . a_{j+1} \cdots a_m \overline{b_1 \cdots b_n}$$

for some integers $m, n, j$, with $m \geqslant 0$, $n \geqslant 1$, and $0 \leqslant j \leqslant m$.
We consider the decimal integers

$$A = a_0 \cdots a_m = \sum_{k=0}^{m} a_k 10^{m-k} \qquad B = b_1 \cdots b_n = \sum_{k=1}^{n} b_k 10^{n-k}.$$

3. We compute the value of $x$ explicitly in terms of $A$ and $B$. First, we get rid of the aperiodic part; we shift it to the left of the decimal point by multiplying by a suitable power of 10, and then we subtract it off, to get

$$10^{m-j} x - A = 0.\overline{b_1 \cdots b_n}. \tag{6.10}$$

4. Next we compute, using (6.9)

$$
\begin{aligned}
0.\overline{b_1 \cdots b_n} &= \frac{B}{10^n} + \frac{B}{10^{2n}} + \frac{B}{10^{3n}} + \cdots \\
&= B\left(\frac{1}{10^n} + \frac{1}{10^{2n}} + \frac{1}{10^{3n}} + \cdots\right) \\
&= B\sum_{k=1}^{\infty}\left(\frac{1}{10^n}\right)^k = \frac{B}{10^n - 1}.
\end{aligned}
$$

5. In the last step we have used the formula of the sum of the geometric series

$$
\sum_{k=1}^{\infty} q^k = \frac{q}{1-q} \qquad |q| < 1.
$$

From equation (6.10) and the result above, we obtain

$$
10^{m-j}x - A = \frac{B}{10^n - 1}
$$

and hence

$$
x = 10^{j-m}\left(A + \frac{B}{10^n - 1}\right).
$$

The right-hand side of this equation consists of sum and products of rational numbers, and is therefore rational. $\quad\square$

   We examine the main steps in the proof:

1. Since we don't begin the proof straight away, we say so. The key notation (the over-bar) is established first in words, then in symbols. The under-brace highlights the structure of the expression.

2. In the theorem there is a hidden quantifier, and the proof begins accordingly. Then we introduce the notation for $x$, and we say why.

3. At every opportunity, we declare our intentions.

4. In this passage, one must decide what constitutes an appropriate amount of details: we have been conservative.

5. A concise reminder of a well-known summation formula.

   We give a concise version of the same proof.

CONCISE PROOF. Let $x$ be a real number with eventually repeating decimal digits $b_k$. Without loss of generality, we may assume that the integer part of $x$ is zero, and that the fractional part is purely periodic:

$$
x = 0.\overline{b_1 \cdots b_n}. \tag{1}
$$

(Any real number may be reduced to this form by first multiplying by a power of 10, and then subtracting an integer, and neither operation affects the property of having repeated digits.) Defining the decimal integer $B = b_1 \cdots b_n$, we find, from (1)

$$
\begin{aligned}
x &= \frac{B}{10^n} + \frac{B}{10^{2n}} + \frac{B}{10^{3n}} + \cdots = B \sum_{k=1}^{\infty} \left( \frac{1}{10^n} \right)^k \\
&= \frac{B}{10^n - 1}.
\end{aligned}
$$

We see that $x$ is rational. $\square$

The expression 'without loss of generality' indicates that the restriction being introduced does not weaken the argument in any way (see also section 4.2). A brief parenthetic remark is added for clarification: it too could be omitted.

## Exercises

**Exercise 6.1.** The following statements are disguised implications: write them explicitly as implications, in the form

*For all ...x, if x ..., then ...*

introducing an appropriate symbol(s) $x$. For example, the statement

*Every differentiable real function is continuous'*

should be rewritten as

*For all real functions $f$, if $f$ is differentiable, then $f$ is continuous.*

1. *Between any two distinct rationals there is another rational.*

2. *The reciprocal of a positive reduced fraction is reduced.*

3. *The sum of two odd integers is even.*

4. *The sum of two odd functions is an odd function.*

5. *The inverse of an invertible matrix is invertible.*

6. *Every subset of a finite set is finite.*

7. *A bounded real function cannot be surjective.*

8. *No polynomial function of positive degree is bounded.*

9. *Three consecutive odd integers greater than 3 cannot all be prime.*

10. *The composition of two surjective functions is surjective.*

11. *The only regular convex polyhedra are the five Platonic solids.*

**Exercise 6.2.** Each of the following statements is an implication, which may be true or false. For each implication

    *i*)  state whether it's true or false;

   *ii*)  state the contrapositive;

  *iii*)  state the converse, and whether it's true or false;

   *iv*)  state the negation.

If the statement in *i*) and *iii*) is false, give a counterexample.

1. *Two integers which are co-prime are also prime.*

2. *If integer divides the product of two integers, then it divides one of the factors.*

3. *If the derivative of a function is increasing, then the function is increasing.*[2]

4. *If an integer is the product of two primes, then it has four distinct divisors.*[3]

**Exercise 6.3.** You are given cryptic proofs of mathematical statements. Rewrite them in a style appropriate for beginning university students, who may be unfamiliar with proofs, and may have forgotten basic facts. Make sure that the strategy of the proof is clear, that each equation is supported by appropriate comments, and that every detail is carefully explained.

  (a)  *Prove that the line through the point $(4,5,1) \in \mathbb{R}^3$ parallel to the vector $\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$ and the line through the point $(5,-4,0)$ parallel to the vector $\begin{pmatrix} 2 \\ -3 \\ 1 \end{pmatrix}$ intersect at the point $(1,2,-2)$.*

    BAD PROOF.

$$\mathbf{v} = \begin{pmatrix} 4 \\ 5 \\ 1 \end{pmatrix} + \lambda \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 4+\lambda \\ 5+\lambda \\ 1+\lambda \end{pmatrix}$$

$$= \begin{pmatrix} 5 \\ -4 \\ 0 \end{pmatrix} + \mu \begin{pmatrix} 2 \\ -3 \\ 1 \end{pmatrix} = \begin{pmatrix} 5+2\mu \\ -4-3\mu \\ \mu \end{pmatrix};$$

$$4+\lambda = 5+2\mu; \qquad 5+\lambda = -4-3\mu; \qquad 1+\lambda = \mu;$$

---

[2]Increasing here is intended in the strict sense.

[3]You may assume that all relevant integers are positive.

$$4 + \lambda = 5 + 2(1 - \lambda) \Rightarrow \lambda = -3, \mu = -2;$$

this gives $\mathbf{v} = \begin{pmatrix} 1 \\ 2 \\ -2 \end{pmatrix}$. $\square$

[*Introduce symbols for the position vectors of the points of the two lines.*]

(b) *Prove that the real function $x \mapsto 3x^4 + 4x^3 + 6x^2 + 1$ is positive.*

BAD PROOF.

$f'(x) = 12x(x^2 + x + 1) = 0 \Leftrightarrow x = 0;$
$f''(0) > 0$: minimum. $f(0) = 1 > 0$. $\square$

[*Define every symbol you use. The second line requires some care.*]

(c) *Prove that for all real values of $a$, the line*

$$y = ax - \left( \frac{a-1}{2} \right)^2$$

*is tangent to the parabola $y = x^2 + x$.*

BAD PROOF.

$ax - (a-1)^2/4 = x^2 + x;$
$x^2 + x(1-a) + (a-1)^2/4 = 0;$
$x = (a-1)/2.$
For this $x$, $dy/dx$ is the same. $\square$

[*What does it mean for a line to be tangent to a curve? Set an appropriate notation for line and parabola.*]

**Exercise 6.4.** Prove that, for any integer $n > 1$, the integer $1 + n!$ is followed by $n - 1$ composite integers. (This shows that there are arbitrarily large gaps between primes.)

**Exercise 6.5.** Prove that the set operators of union and intersections are associative.

**Exercise 6.6.** Write the first few sentences of the proof of each statements, setting the notation, and identifying the RTP. (For this task, understanding the meaning of the statements is unimportant.)

1. *The order of a finite group is divisible by the order of each one of its subgroups.*

2. *Every valuation of a field with prime characteristic is non-archimedean.*

3. *A bi-infinite sequence over the alphabet $\{0, 1\}$ is Sturmian if and only if it is balanced and not eventually periodic.*

4. *The characteristic polynomial of the incidence matrix of a Pisot substitution is irreducible over $\mathbb{Q}$.*

5. *Any hyperbolic matrix in $SL_2(\mathbb{Z})$ is conjugate to a matrix with non-negative coefficients.*

6. *Any irrational number can be expressed in just one way as an infinite simple continued fraction.*

7. *A subgroup is normal if it contains all the conjugates of each one of its elements.*

8. *A subset of a metric space is open if and only if its complement is closed.*

9. *On a compact set every continuous function is uniformly continuous.*

# Chapter 7

# Induction

Induction is one of the most basic methods of proof, used in every area of mathematics. It concerns statements which can be formulated as follows

*All natural numbers have property $\mathscr{P}$.*

Usually, induction is first met in the proof of summation identities and inequalities, such as

1) $\displaystyle\sum_{k=0}^{n} x^k = \frac{1 - x^{n+1}}{1 - x}$ $\qquad x \neq 1$ $\qquad$ sum of geometric progression

2) $\displaystyle (x+y)^n = \sum_{k=0}^{n} \binom{n}{k} x^k y^{n-k}$ $\qquad\qquad$ binomial theorem

3) $(1+x)^n \geqslant 1 + nx$ $\qquad\qquad x > -1$ $\qquad$ Bernoulli inequality

4) $\displaystyle\prod_{k=1}^{n} x_k \leqslant \left(\frac{1}{n} \sum_{k=1}^{n} x_k\right)^n$ $\qquad\qquad x_k > 0$ $\qquad$ AGM inequality

5) $\displaystyle\left|\sum_{k=1}^{n} x_k y_k\right|^2 \leqslant \sum_{k=1}^{n} |x_k|^2 \sum_{k=1}^{n} |y_k|^2$ $\qquad$ Cauchy-Schwarz inequality

The variables are complex numbers in items 1,2, and real numbers in 3–5. (AGM stands for *arithmetico-geometric mean*.)

By encoding property $\mathscr{P}$ as a predicate over $\mathbb{N}$, we can write this statement symbolically as

$$\forall n \in \mathbb{N}, \ \mathscr{P}(n). \tag{7.1}$$

This expression unfolds into an infinite sequence of boolean expressions

$$\mathscr{P}(1), \mathscr{P}(2), \mathscr{P}(3), \ldots.$$

To prove it, we must show that all expressions in the sequence are true:

$$T, T, T, \ldots.$$

For instance, the Bernoulli inequality unfolds as follows

$$1+x \geqslant 1+x, \quad (1+x)^2 \geqslant 1+2x, \quad (1+x)^3 \geqslant 1+3x, \quad \ldots$$

each inequality being valid for all real numbers $x > -1$.

A special form of mathematical argument is needed for this purpose, called **mathematical induction**. It takes several forms, and we will examine four:

(*A*) The well-ordering (or 'least counterexample') principle.

(*B*) The infinite descent method.

(*C*) The induction principle.

(*D*) The strong induction principle.

We will show that these principles are equivalent, in the sense that any one of them implies all the others. The different forms of the principle are appropriate for different contexts.

The term mathematical induction was coined by Augustus De Morgan (around 1840). The adjective 'mathematical' is used to differentiate this concept from the homonymous concept in philosophy, which has quite a different meaning[1]. The induction principle (in any of the above forms) is one of the five **Peano axioms**[2] which characterise the structure of the set $\mathbb{N}$ of natural numbers.

---

[1]An inductive argument in philosophy is where one uses very strong premises to support the *probable* truth of a conclusion.

[2]Giuseppe Peano, Italian mathematician (1858–1932).

# 7.1   The well-ordering principle

The **well-ordering principle** states that:

(A) *Every non-empty set of natural numbers contains a least element.*

This principle is clearly false for subsets of $\mathbb{Z}, \mathbb{Q}$ or $\mathbb{R}$, and, less trivially so, for subsets of $\mathbb{Q}^+$ or $\mathbb{R}^+$, or of any real or rational interval. A possible strategy for proving a statement of the form (7.1) is to combine well-ordering with contradiction. A proof of this kind will begin as follows:

PROOF.   Suppose (7.1) is false. Let $k$ be the least natural number $n$ for which $\mathscr{P}(n)$ is false. ...

Then, with this $k$, we try to deduce a contradiction, using the knowledge that $\mathscr{P}(n)$ is true for $n = 1, \ldots, k-1$. The integer $k$ defined in the proof exists by virtue of the well-ordering principle: it is the least element of the set $\{n \in \mathbb{N} : \neg\mathscr{P}(n)\}$, which is a subset of $\mathbb{N}$. By assumption, this subset is non-empty.

Let us apply this strategy to some specific problems.

*Show that if $n \geqslant 7$, then $n! > 3^n$.*

PROOF.   Suppose this statement is false, and let $k$ be the least integer $n \geqslant 7$ such that $n! \leqslant 3^n$. We know that $k > 7$, since

$$7! = 5040 > 2187 = 3^7.$$

Put $j = k-1$. Then $j \geqslant 7$, and hence $j! > 3^j$ by choice of $k$. So

$$k! = (j+1)j! > (j+1)3^j > 3(3^j) = 3^k$$

contradicting the choice of $k$.   □

The initial verification that $k > 7$ creates enough room for the principle to work. We note that $6! = 720$ and $3^6 = 729$, so the inequality $n \geqslant 7$ in the statement above is the best possible one. Under these circumstances, we say that $n \geqslant 7$ is a **strict bound**.

**Theorem.**   *Every integer greater than 1 is a product of one or more prime numbers.*

PROOF.   Suppose the statement is false, and let $k$ be the smallest integer greater than 1 which is not a product of primes. Then $k$ is not prime, so $k = mn$ for two smaller integers $m, n \geqslant 2$. Since $m$, $n$ are smaller than $k$ and greater than 1, they are

products of prime numbers. So $k$ is a product of primes too, because $k = mn$. This contradicts the choice of $k$. $\quad\square$

This form of mathematical induction can be very effective. The assumption that $k$ is a smallest counterexample puts the maximum amount of information on the table; it gives us the feel that we are examining a concrete object.

## 7.2   The infinite descent method

This is a second form of mathematical induction, developed in the 17th Century by Fermat[3]. It exploits the following variant of the well-ordering principle:

(B) *There is no infinite decreasing sequence of natural numbers.*

This principle follows from the well-ordering principle $((A) \Rightarrow (B))$: if there were an infinite decreasing sequence of natural numbers

$$n_1 > n_2 > n_3 > \cdots$$

then the set

$$\{n_1, n_2, n_3, \ldots\},$$

which is a non-empty subset of $\mathbb{N}$, would have no smallest element. Conversely, well-ordering follows from descent $((B) \Rightarrow (A))$ —see exercises.

As with well-ordering, we combine descent with contradiction. To prove that all natural numbers have a certain property $\mathscr{P}$, we suppose that some natural number $n_1$ doesn't have property $\mathscr{P}$, and we try to show that some smaller natural number $n_2$ also doesn't have property $\mathscr{P}$. If we succeed, then, repeating this procedure for $n_2$, etc., we obtain a descending sequence of natural numbers $n_1 > n_2 > n_3 > \cdots$ none of which have property $\mathscr{P}$. But then the set $\{n \in \mathbb{N} : \neg\mathscr{P}(n)\}$ has no least element, contradicting well-ordering.

Fermat used extensively this method to show that certain properties or relations are impossible for whole numbers. Infinite descent was brought to prominence by Euler, who proved with it that it is impossible to find natural numbers $x, y, x$ such that

$$x^3 + y^3 = z^3.$$

This is a special case of the celebrated **Fermat's last theorem** (where the exponent 3 is replaced by any integer greater than 2).   In keeping with this tradition, let us

---

[3]Pierre de Fermat, French lawyer and amateur mathematician (1601 or 1607/8–1665)

apply the method of infinite descent to prove the non-existence of integer solutions of a polynomial equation.

**Theorem.** *There are no natural numbers $m, n$ such that $2n^2 = m^2$.*

This statement is equivalent to the irrationality of $\sqrt{2}$ (see section 6.1), as one sees by dividing each term of the equation by $n^2$, and then taking the positive square root.

PROOF. Suppose that there are natural numbers $n_1, m$, such that $2n_1^2 = m^2$. So 2 divides $m^2$, and hence 2 divides $m$, so that $m = 2n_2$ for some natural number $n_2$. Then $2n_1^2 = m^2 = 4n_2^2$, and hence $n_1^2 = 2n_2^2$. So $n_2 < n_1$ and $2n_2^2$ is the square of a natural number.

Repeating the argument, there is $n_3 < n_2$ such that $2n_3^2$ is the square of a natural number. This continues for ever, which is impossible. Hence the equation $2n^2 = m^2$, has no solution, as stated. $\square$

# 7.3   Peano's induction principle

This is a third form of mathematical induction; it was proposed by Dedekind[4] and formalised by Peano. The **principle of induction** states that:

(C) *If $\mathscr{P}$ is a predicate over $\mathbb{N}$ such that*

    *i)  $\mathscr{P}(1)$ is true;*
    *ii)  for all $k \in \mathbb{N}$, $\mathscr{P}(k) \Rightarrow \mathscr{P}(k+1)$*

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

The first condition is called the **base case**;  the second the **inductive step.** The bi-unique correspondence between predicates over $\mathbb{N}$ and subsets of $\mathbb{N}$, given by

$$S = \{n \in \mathbb{N} : \mathscr{P}(n)\}, \qquad \mathscr{P}(x) = (x \in S),$$

allows us to reformulate the induction principle in the language of sets.

(C′) *Let $S$ be a subset of $\mathbb{N}$. If*

    *i)   $1 \in S$        ii)   $\forall k \in \mathbb{N}, \ k \in S \Rightarrow (k+1) \in S$*

  *then $S = \mathbb{N}$.*

---

[4]Richard Dedekind, German mathematician (1831–1916).

The principle of induction follows from well-ordering $((A) \Rightarrow (C))$. To show this, we define the set

$$S' = \mathbb{N} \setminus S = \{n \in \mathbb{N} : \neg \mathscr{P}(n)\}.$$

We prove that $S'$ is empty by contradiction. Assume that $S'$ is non-empty. By the well-ordering axiom, $S'$ has a least element $m$. Since $1 \in S$, then $m \neq 1$, and since $m$ is the smallest element of $S'$, it follows that $m - 1 \in S$. Then, putting $k = m - 1$, we find from condition $ii$) that $k + 1 = m \in S$, which contradicts the fact that $m \in S'$.

The base case for Peano's axiom could be any integer, not just 1. Indeed if $\mathscr{P}(n)$ is valid for all $n \geqslant k$, then by letting $m = n - k + 1$, the range $n \geqslant k$ becomes $m \geqslant 1$.

Peano's induction provides straightforward proofs of finite sums and products formulae:

$$\forall n \in \mathbb{N}, \qquad \sum_{k=1}^{n} a_k = F(n) \tag{7.2}$$

$$\forall n \in \mathbb{N}, \qquad \prod_{k=1}^{n} a_k = G(n) \tag{7.3}$$

where $(a_k)$ is a sequence, and $F$ and $G$ are explicit functions of $n$, hopefully easier to compute than the original sum or product.

In an inductive proof of (7.2), the base case consists of verifying that $a_1 = F(1)$. In the inductive step, assuming that the formula holds for some $n$, we obtain

$$\sum_{k=1}^{n+1} a_k = \sum_{k=1}^{n} a_k + a_{n+1} = F(n) + a_{n+1}.$$

Thus a proof by induction of (7.2) reduces to the proof of two statements

$$a_1 = F(1), \qquad F(n) + a_{n+1} = F(n+1) \quad n \geq 1 \tag{7.4}$$

which no longer involves summation. For example, to prove the formula

$$\sum_{k=1}^{n} k^2 = \frac{n(n+1)(2n+1)}{6} \qquad n = 1, 2, \ldots \tag{7.5}$$

we must verify the following statements

$$1^1 = \frac{1 \cdot 2 \cdot 3}{6} \qquad \frac{n(n+1)(2n+1)}{6} + (n+1)^2 = \frac{(n+1)(n+2)(2n+3)}{6}.$$

For a product formula, the expressions (7.4) are replaced by

$$a_1 = G(1), \qquad G(n) \times a_{n+1} = G(n+1) \quad n \geq 1.$$

These identities can be verified using a computer algebra system.

A similar arrangement applies to the inductive proof of inequalities. However, even in the simplest cases, such proofs are not as mechanical as those of summation and product formulae.

*Prove that, if $n > 3$, then $2^n \geq n^2$.*

PROOF. We prove it by induction on $n$. The base case $n = 4$ is readily verified: $2^4 \geq 4^2$. Assume now that for some $k \geq 4$ we have $2^k \geq k^2$. Then

$$2^{k+1} = 2 \cdot 2^k \geq 2 \cdot k^2$$

where the inequality follows from the induction hypothesis. To complete the proof, we must show that, in the $k$-range under consideration, we have $2 \cdot k^2 \geq (k+1)^2$. Indeed, the polynomial

$$p(k) = 2 \cdot k^2 - (k+1)^2 = k^2 - 2k - 1$$

has roots $1 \pm \sqrt{2}$, with

$$-1 < 1 - \sqrt{2} < 1 + \sqrt{2} < 3.$$

Thus, for $k \geq 3$, we have $p(k) > 0$, or $2k^2 > (k+1)^2$. Hence $2^{k+1} \geq (k+1)^2$, completing the induction. $\square$

A straightforward inductive strategy doesn't always work, as the following example shows.

*Prove that, for all integers $n \geq 0$ and real numbers $x > -1$, we have*

$$(1+x)^n > nx. \tag{7.6}$$

We try to prove this by induction on $n$.

Base case: $n = 0$. We have $(1+x)^0 = 1 > 0 = 0x$, as desired.
Inductive step: Assume (7.6) for some $n = k \geq 0$. Then

$$\begin{aligned}
(1+x)^{k+1} &= (1+x)(1+x)^k \\
&> (1+x)kx \quad \text{(by induction hypothesis)} \\
&= kx + kx^2 \quad (??)
\end{aligned}$$

We don't seem to be getting anywhere.

PROOF.   We prove the *stronger* inequality $(1+x)^n \geqslant 1+nx$. (This is the **Bernoulli inequality**, see page 148.)
Base case: $n = 0$. Since $x > -1$, we have $(1+x)^0 = 1 \geqslant 1+0x$.
Inductive step: Assume it for $n = k \geq 0$. Then

$$\begin{aligned}
(1+x)^{k+1} &= (1+x)(1+x)^k \\
&\geqslant (1+x)(1+kx) \quad \text{(by induction hypothesis)} \\
&= 1+(k+1)x+kx^2 \geqslant 1+(k+1)x.
\end{aligned}$$

The inductive step is complete.   □

## 7.4   Strong induction

The conversion from well-ordering to induction illustrated in the previous section doesn't always work. To see why, let us try to prove by the induction principle that every integer greater than 1 is a product of one or more prime numbers.

This statement is true for 2, and for 3. How does 3 being prime help us show that 4 is a product of prime numbers? Not at all! The rescue is a form of induction called **strong induction** (or **complete induction**), which is formulated as follows

(D) *Let $\mathscr{P}$ be a predicate on $\mathbb{N}$. If*

  i) $\mathscr{P}(1)$ *is true;*
  ii) *for any $k \in \mathbb{N}$, if $\mathscr{P}(1), \mathscr{P}(2),\ldots,\mathscr{P}(k)$ are true, so is $\mathscr{P}(k+1)$;*

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

As with Peano's induction principle, the base case for strong induction can be any integer. Furthermore, the principle may be reformulated in terms of sets

(D′) *Let $S$ be a subset of $\mathbb{N}$. If*

  i) $1 \in S$;
  ii) *for any $k \in \mathbb{N}$, if $\{1,2,\ldots,k\} \subset S) \Rightarrow (k+1) \in S$*

  *then $S = \mathbb{N}$.*

Strong induction follows from Peano's induction $((C) \Rightarrow (D))$, since the latter has fewer conditions. (One could say that induction is *stronger* than strong induction!)

**Theorem.** *Every integer greater than 1 is a product of one or more prime numbers.*

PROOF. We use strong induction. The integer 2 is a prime, which establishes the base case. Suppose that $n$ is an integer greater than 1, and let us assume that every integer $k$ with $2 \leqslant k < n$ is a product of primes. We have two cases:

*Case I:* $n$ is a prime. Then $n$ is a product of primes.
*Case II:* $n$ is not a prime. Then there are $x, y > 1$ with $xy = n$. By induction hypothesis, $x$ and $y$ are products of primes, and hence so is $n$. □

To complete the proof that the four formulations of the induction principle are equivalent, it remains to show that well-ordering follows from induction $((D) \Rightarrow (A))$. Indeed, after proving this implication, we will have proved that

$$(A) \Leftrightarrow (B) \qquad (A) \Rightarrow (C) \Rightarrow (D) \Rightarrow (A).$$

Equivalence is thus established by means of a circular argument (cf. section 6.5).

We use the both ends method. Assume that strong induction holds but that there is a non-empty subset $S$ of $\mathbb{N}$ without least element. Let $S' = \mathbb{N} \smallsetminus S$. Then $1 \in S'$, for otherwise 1 would be the least element of $S$. Likewise, if, for some $k \geqslant 1$ we have $\{1, \ldots, k\} \subset S'$, then we must also have $k + 1 \in S'$, because otherwise $k + 1$ would be the smallest element of $S$. Hence, from strong induction, we have that $S' = \mathbb{N}$, and hence $S = \emptyset$, contradicting the hypothesis $S \neq \emptyset$.

## 7.5 Good manners with induction proofs

Induction proofs are easy to structure, as there is a clear procedure to follow. Even if the various steps comprising such a proof are predictable, it's important that they be spelled out clearly.

At the beginning of the proof, say with respect to what variable induction is performed.

*We proceed by induction on the degree $n$ of the polynomial.*

Say at what point you are using the induction hypothesis.

*... where the last inequality follows from the induction hypothesis.*

If an induction argument is easy, one may be tempted to skip it altogether. This should be done with some care, as one must judge the mathematical maturity of the audience.

Let $x_1 = 2$, and let $x_{n+1} = x_n^2$, for $n \geq 1$. Then $x_n = 2^{2^{n-1}}$.

The last sentence is a bit blunt. It should be replaced by more considerate sentences such as

A straightforward induction argument shows that $x_n = 2^{2^{n-1}}$.
Then $x_2 = 2^2$, $x_3 = 2^{2^2}$, and, in general, $x_n = 2^{2^{n-1}}$.

When using the Peano form (or strong induction) on a variable $n$, it is considered good practice to go from $n = k$ to $n = k+1$ rather than from $n$ to $n+1$. This arrangement removes any ambiguity. Since a new symbol is also introduced, this notation may not be appropriate all cases.

## Exercises

**Exercise 7.1.**

1. State and prove Peano's induction principle for an arbitrary base case, not necessarily unity.

2. Prove that the well-ordering principle follows from strong induction.

**Exercise 7.2.** or a natural number $n$, let $\mathcal{L}$ be a boolean expression of the type

$$\mathcal{L} = \pi_1 x_1 \in X_1, \pi_2 x_2 \in X_2, \cdots, \pi_n x_n \in X_n, \ \mathcal{P}(x_1, \ldots, x_n)$$

where the $\pi_k$ are quantifiers ($\pi_k \in \{\forall, \exists\}$), the $X_k$ are sets, and $\mathcal{P}$ is a predicate on $X_1 \times \cdots \times X_n$. Prove, by induction on $n$, that $\neg \mathcal{L}$ is obtained from $\mathcal{L}$ by interchanging quantifiers and negating $\mathcal{P}$.

**Exercise 7.3.** Prove by induction the following summation formulae:

$$\sum_{k=1}^{n} (2k - 1) = n^2$$

$$\sum_{k=1}^{n} k^3 = \left( \sum_{k=1}^{n} k \right)^2$$

$$\sum_{k=1}^{n} k! \cdot k = (n+1)! - 1$$

$$\sum_{k=1}^{n} \frac{1}{k^2 - 1} = \frac{3}{4} - \frac{2n+1}{2n(n+1)}.$$

**Exercise 7.4.** Prove by induction.

1. The sum of arithmetic and geometric progressions.

2. The binomial theorem.

3. The AGM inequality.

4. The Cauchy-Schwarz inequality.

5. Leibniz theorem for the higher-order derivatives of the product of two functions.

6. The chain rule for differentiation of the repeated composition of functions.

# Chapter 8

# Existence and definitions

An existence statement asserts that there exists a quantity $x$ that has a certain property $\mathscr{P}$. In symbols

$$\exists x \in X, \ \ \mathscr{P}(x)$$

where $X$ is a set and $\mathscr{P}$ is a predicate over $X$ (see section 3.4). For example, the statement

*The integer $10^9$ is the sum of two primes*

asserts the existence of a pair of primes with a certain property (what are $X$ and $\mathscr{P}$ in this case?).

In this chapter we discuss proofs of existence.

## 8.1   Proofs of existence

To prove that something exists, it is not necessary that the object in question be constructed explicitly. For example, Euclid's proof of the infinitude of primes (page 133) establishes the existence of infinitely many things, without producing any of them! By contrast, equation (6.8), which proves that Euler's polynomial $n^2 + n + 41$ is composite for infinitely many values of the indeterminate $n$, exhibits such values explicitly.

Accordingly, existence proofs are said to be **effective** (or **constructive**) if an explicit construction is given, and **non-effective** (or **non-constructive**) if no explicit construction is given. The effective method may be described as a two-stage process:

(WHAT?)    Identify an element $x$ of $X$.

(WHY?)     Show that $\mathscr{P}(x)$ is TRUE.

The format of the proof should make clear which is the WHAT part and which is the WHY.

By contrast, all kinds of arguments are used in non-effective existence proofs. Many (but not all!) non-constructive proofs employ contradiction: one assumes that the object being defined does not exist, and derive a false statement.

We begin with is an effective existence proof:

**Theorem.**    *There is a real number $\alpha$ such that $\alpha^2 = 2$.*

PROOF (Liebeck).   Draw a square of side 1.

(WHAT?)    Let $\alpha$ be the length of a diagonal of the square.

(WHY?)     Then by Pythagoras, $\alpha^2 = 2$.       $\square$

This minimalist proof takes for granted that the length of a diagonal of the unit square is a real number. A justification of this step requires the tools of analysis.

Let us re-visit equation (6.8) concerning Euler's polynomial.

**Theorem.**    *There are infinitely many integers $n$ for which $p(n) = n^2 + n + 41$ is composite.*

PROOF.

(WHAT?)    Let $n = 41k$, for $k \in \mathbb{N}$.

(WHY?)     Then $p(41k) = 41(41k^2 + k + 1)$.      $\square$

In this proof the WHAT part identifies an infinite set of integers.

The next example is a non-effective existence proof which does not use contradiction.

**Theorem.**    *There is an irrational number $a$ such that $a^{\sqrt{2}}$ is rational.*

PROOF.  Consider $\sqrt{2}^{\sqrt{2}}$. We have two cases:

CASE I:   $\sqrt{2}^{\sqrt{2}}$ is rational. Then $a = \sqrt{2}$ is as required.

CASE II:  $\sqrt{2}^{\sqrt{2}}$ is irrational. Then

$$\left( \sqrt{2}^{\sqrt{2}} \right)^{\sqrt{2}} = \sqrt{2}^{\sqrt{2}\sqrt{2}} = \sqrt{2}^2 = 2,$$

so $a = \sqrt{2}^{\sqrt{2}}$ is as required.      $\square$

This is a proof by cases —see section 6.2. We note that in this proof we aren't given a WHAT. There are two different WHATS, and we have no idea which of them works for the WHY.

Establishing the existence of solutions of an equation is a common task. When we attempt to solve an equation, often we are looking for a particular number, in which case we are after an effective existence proof. In other cases, we may be satisfied by a proof that a solution exists at all, or exists in a specified range of values of the argument. We now present an effective and a non-effective existence proof of the same statement.

**Theorem.** *The equation $x^4 - 10x^2 + 1 = 0$ has four distinct real solutions.*

EFFECTIVE EXISTENCE PROOF. Let $f(x) = x^4 - 10x^2 + 1$.
(WHAT?)   Let $x_\pm = \sqrt{3} \pm \sqrt{2}$; then $x_\pm$ are real numbers.
(WHY?)   We compute, using the binomial theorem

$$
\begin{aligned}
f(x_\pm) &= (\sqrt{3} \pm \sqrt{2})^4 - 10(\sqrt{3} \pm \sqrt{2})^2 + 1 \\
&= 9 \pm 12\sqrt{6} + 36 \pm 8\sqrt{6} + 4 - 30 \mp 20\sqrt{6} - 20 + 1 \\
&= 0
\end{aligned}
$$

where the top and bottom signs in $\pm$ and $\mp$ match. The above calculation shows that $x_\pm$ are roots of $f$. Since the real numbers $x_\pm$ are positive and distinct, and $f$ is an even function, we conclude that $-x_\pm$ are also roots of $f$.   □

This effective existence proof involved 'guessing' the solutions, and then verifying that the guess was correct. Mathematicians often hide from view the reasons behind a guess, but I wouldn't encourage this practice. Here there was no 'guess'. I worked backwards, starting from $x_+$, and then using linear algebra to construct the quartic equation of which $x_+$ is a solution.

NON-EFFECTIVE EXISTENCE PROOF. Let $f(x) = x^4 - 10x^2 + 1$. The computation

$$ f(0) = 1 \quad f(1) = -8; \quad f(4) = 96, $$

shows that $f$ changes sign twice in the interval $(0,4)$. Because $f$ is continuous, there exist two distinct real numbers $x_-, x_+$ in the open intervals $(0,1)$ and $(1,4)$, respectively, at which the function $f$ vanishes. The real numbers $-x_\pm$ are also roots of $f$, because $f$ is even. The four numbers $\pm x_\pm$ are clearly distinct.   □

The key argument in the proof was inferring the existence of a root from a change of sign of a continuous function. This is the **intermediate value theorem**

from real analysis. The proof provides some information about the solutions, in the form of bounds. These bounds can be sharpened by doing some extra work, e.g.,

$$f\left(\frac{7}{22}\right) = \frac{-503}{234256} \qquad f\left(\frac{6}{19}\right) = \frac{1657}{130321}$$

and hence

$$\frac{6}{19} < x_- < \frac{7}{22}.$$

Considering that $7/22 - 6/19 \approx 0.002$, our knowledge of the solution $x_-$ has increased markedly. We see that the distinction between an effective and a non-effective proof is blurred.

**Theorem.** *If two integers are the sum of two squares, then so is their product.*

This statement is an implication. The existence statement in the deduction is conditional to another existence statement in the hypothesis. We give a concise effective proof, where an eloquent formula supplies at once the WHATS and the WHYS.

PROOF. This statement follows from the algebraic identity

$$(j^2 + k^2)(m^2 + n^2) = (jm + kn)^2 + (km - jn)^2$$

which is established by a straightforward calculation. $\square$

Our next example is a well-known existence statement, **Dirichlet's box principle** (or the **pigeon-hole principle**).

**Theorem.** *Given $n$ boxes and $m$ objects in them, if $m > n$, then at least one box contains more than one object.*

Even though this statement is obvious, we give a proof, necessarily non-effective. We use contradiction, namely the **both ends method** described in section 6.6. It's easy to predict where contradiction will lead us.

PROOF. We use contradiction. Let $m$ and $n$ be integers, and let $m_j$ be the number of objects in the $j$th box. If $m > n$ and $m_j \leqslant 1$, for $j = 1, \ldots, n$, then

$$m = \sum_{j=1}^{n} m_j \leqslant \sum_{j=1}^{n} 1 = n$$

giving $m \leqslant n$, contrary to the assumption. $\square$

In this proof we are not given the WHAT, and hence there cannot be a WHY step either. The proof begins by introducing the symbol $m_j$; this simple but important step makes the rest immediate.

Dirichlet's principle looks informal, as it deals with 'objects in boxes', or even 'pigeon-holes'. It is, however, a crisp mathematical statement, which may be interpreted geometrically. Let us consider the $n$-dimensional integer vectors

$$\mathbf{m} = (m_1, \ldots, m_n) \in \mathbb{N}^n$$

where we stipulate that the set $\mathbb{N}$ contains zero. We consider the following subset of $\mathbb{N}^n$

$$\Omega_m := \{(m_1, \ldots, m_n) \in \mathbb{N}^n : \sum_j m_j = m\}.$$

Dirichlet's principle says that if $m > n$, then every point $\mathbf{m}$ in $\Omega_m$ lies outside the $n$-dimensional unit cube.

We consider an application of Dirichlet's box principle.

**Theorem.** *Given a positive integer $n$, in any set of integers with more than $n$ elements, there must be two integers whose difference is divisible by $n$.*

PROOF. Let $n$ be given, and let $S$ be set of $m$ integers, with $n < m$. We divide each element of $S$ by $n$, obtaining $m$ integer remainders. These remainders can assume at most $n$ distinct values, and therefore two elements of $S$, say $x_i$ and $x_j$ must have the same remainder, by Dirichlet's box principle. But then $x_i - x_j$ gives reminder zero when divided by $n$, as desired. $\square$

## 8.2 Unique existence

A statement of unique existence has the form

> *There is exactly one element of $X$ having property $\mathscr{P}$.*

where, as usual, $X$ is a set, and $\mathscr{P}$ is a predicate over $X$. For example

> *The equation $f(x) = 0$ has exactly one solution.*
> *The sets $A$ and $B$ have precisely one point in common.*
> *The function $g$ has a unique stationary point.*

The adverbs 'exactly' and 'precisely' differentiate unique existence from existence.

Unique existence is a hidden conjunction. The conjuncts are:

EXISTENCE:   *There is one element of $X$ with property $\mathscr{P}$.*
UNIQUENESS:  *If $x, y \in X$ have property $\mathscr{P}$, then $x = y$.*

Expressing unique existence with symbols is cumbersome:

$$[\exists x \in X, \ \mathscr{P}(x)] \wedge [\forall x, y \in X, \ (\mathscr{P}(x) \wedge \mathscr{P}(y)) \Rightarrow (x = y)] \qquad (8.1)$$

and for this reason, the special symbol $\exists!$ is used for this purpose

$$\exists! x \in X, \ \mathscr{P}(x).$$

In a proof of unique existence, the conjuncts should normally be proved separately, in either order.

**Theorem.**   *The identity element of a group is unique.*

The group axioms were given in section 6.1. The existence of an identity element does not require a proof, being an axiom (axiom G3). We only have to prove uniqueness.

PROOF.  Let us assume that a group has two identity elements, $\Diamond_1$ and $\Diamond_2$. Applying the axiom G3 to each identity element, we obtain

$$\Diamond_1 \odot \Diamond_2 = \Diamond_1 \qquad \Diamond_1 \odot \Diamond_2 = \Diamond_2$$

and hence $\Diamond_1 = \Diamond_2$.   $\square$

**Theorem.**   *For every set $X$ and subset $Y$, there is a unique set $Z$ such that $Y \cup Z = X$ and $Y \cap Z = \emptyset$.*

PROOF.  Let $Z = X \smallsetminus Y$. By construction, $Y \cap Z = \emptyset$. Since $Y \cup Z$ is a subset of $X$, and every element of $X$ belongs to either $Z$ or $Y$, we have $Y \cup Z = X$. So a set $Z$ with the required properties exists.

To prove uniqueness, assume that $Z$ is a set with the stated properties. Then, since $Y \cup Z = X$, $Z$ must be a subset of $X$, and it must contain all elements of $X$ that aren't in $Y$. So $Z$ must contain $X \smallsetminus Y$.

Since $Y \cap Z = \emptyset$, the set $Z$ can't contain any elements of $X$ that are in $Y$. So $Z$ must be $X \smallsetminus Y$.   $\square$

The existence part of this proof is constructive. For this reason, the proof of uniqueness does not establish the second conjunct in (8.1) explicitly. Rather that considering two objects with the stated properties, and show that they are the same, we consider a single object, and show that it is the same as that defined in the existence part.

The **fundamental theorem of arithmetic** is a prominent statement of unique existence.

**Theorem.** *Every natural number greater than 1 may be written as a product of prime numbers. This representation is unique, apart from re-arrangements of the factors.*

The expression 'apart from re-arrangements of the factors' gives the impression of lack of uniqueness. This is not the case: the theorem states that there exists a unique **multiset** of prime numbers with the property that the product of its elements is the given natural number. Alternatively, there is a unique **set** of pairwise coprime prime-powers with the same property.

## 8.3 Definitions

Definitions are instances of unique existence. When we define a mathematical object, we must ensure that this object actually exists, and that our definition identifies it uniquely. If this is the case, then the object is **well-defined.**

In the case of sets, there is a safety net which may protect us from a careless definition. If the conditions imposed on the set are too restrictive, then the set will be empty. In this case the object being defined still exists; however, the consequences of an empty definition may be more serious than formally correct nonsense, e.g., if the set in question is the domain of a function.

> *Let M be the set of 2 by 2 integral matrices, with odd entries and unit determinant.*

The determinant of a matrix with odd entries is an even integer, and so it cannot be 1. The matrices with odd entries form an infinite subset of the set of all 2 by 2 integral matrices; so do the matrices with unit determinant. But the intersection $M$ of these two sets is empty!

More generally, a set defined by several conditions could be presented as follows.

> *Let $A_1, \ldots, A_n$ be sets, and let $A$ denote their common intersection*

$$A = \bigcap_{k=1}^{n} A_k = A_1 \cap A_2 \cap \cdots \cap A_n. \tag{8.2}$$

There are some subtleties in this definition. For $n > 2$, the set $A$ is well-defined, because the intersection operator is associative —see exercises. For $n = 1$, formula (8.2) is undefined, although in this case we can reasonably assume that the formula means $A = A_1$. If $n \leqslant 1$ it would be legitimate to interpret the formula as $A = \emptyset$, thereby extending (8.2) to all integers. However, even if $n \geq 1$, the possibility that $A$ be empty is quite evident here. For example, if $A_k$ is the set of solutions of an equation, then $A$ is the set of solutions of a system of $n$ simultaneous equations, and there may be no solution.

Let us write a more considerate definition of the set $A$, which accounts for all eventualities, and alerts the reader to a potential danger.

> Let $A_1, \ldots, A_n$ be sets. Let $A$ be the (possibly empty) common intersection of these sets (with $A = A_1$, if $n = 1$, and $A = \emptyset$, if $n < 1$).

In the next example, we define a sequence whose existence is surprisingly nontrivial.

> Let $p_n$ be the smallest prime number with $n$ decimal digits.

Is this sequence well-defined? For this to be the case, the prime $p_n$ must exist for all $n \in \mathbb{N}$. Let's look at the first few terms:

$$p_1 = 2, \quad p_2 = 11, \quad p_3 = 101, \quad p_4 = 1009, \quad p_5 = 10007, \ldots$$

In all cases, the prime $p_n$ lies just above $10^{n-1}$, and it seems overwhelmingly likely that at least one prime exists between $10^{n-1}$ and $10^n$. On the other hand, we do know that there are arbitrarily large gaps between consecutive primes (see exercise 4). There are large gaps in our sequence, e.g., $p_{21} = 10^{21} + 117$: could it possibly happen that a gap include all integers with $n$ digits, for some $n$?

It turns out that this cannot happen, and our sequence is indeed well-defined. But the proof of this requires some sophisticated results on the size of prime numbers. (An accessible introduction to this beautiful part of mathematics is found in [10]. See also exercise 5, below).

A function definition $f : A \rightarrow B$ requires specifying two sets, the domain $A$ and the co-domain $B$, together with a rule that associates to every point $x \in A$ a unique point of $f(x) \in B$. When we define a function, somebody can always ask us

*How do you know that this function exists?*
*How do you know that its domain is A?*
*How do you know that its co-domain is B?*

The specification of $A$ and of the rule $x \mapsto f(x)$ must not contain any ambiguity, lest $f$ is not a function. By contrast, there is some flexibility in the specification of the co-domain $B$, in the sense that any set containing $f(A)$ may serve as a co-domain. Formally, different choices of the co-domain correspond to different functions; in practice, however, such distinctions are often unimportant.

If the domain of a function $f$ is $A$, why don't we always choose $f(A)$ as co-domain? This would have the advantage of making every function surjective! The problem is that we may not know what $f(A)$ is, or the description of $f(A)$ may be exceedingly complicated.

For instance, consider the function

$$f : \mathbb{N} \to \mathbb{N} \qquad n \mapsto p_n - 10^{n-1} \tag{8.3}$$

where $p_n$ is the sequence of primes given above. We find

| $n$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| $f(n)$ | 1 | 1 | 1 | 9 | 7 | 3 | 3 | 19 | 7 | 7 |

This function is well-defined, but we know little about the image $f(\mathbb{N})$. Clearly $f(\mathbb{N})$ cannot contain any integer $n$ divisible by 2 or by 5 (see exercise 4), but we don't even know if $f(\mathbb{N})$ is bounded.

## 8.4 Recursive definitions

The recursive definition of a sequence is a form of definition which is closely connected to the **principle of induction**. This method can deliver great complexity from minimal ingredients.

We have pointed out that a sequence $(a_k)$ corresponds to a function over $\mathbb{N}$, whereby $a_k$ represents the value of the function at the natural number $k$. However, it is unclear how this function may be extracted from the definition of the sequence. For example, we don't know any 'useful' way to express the $k$th prime number as an explicit function of $k$. In absence of an explicit formula, we seek to represent $a_k$ in terms of $a_{k-1}$. These are the **recursive sequences**, which are specified by two data: the first term of the sequence, called the **initial condition**, and the rule which gives a term from the previous one.

For instance, the **factorial sequence** $n! = 1 \cdot 2 \cdot 3 \cdots n$, could also be defined recursively as follows.

$$0! = 1 \qquad\qquad k! = k \cdot (k-1)! \qquad k \geqslant 1. \qquad\qquad (8.4)$$

This definition, in effect, specifies how the product is computed. Likewise, the $k$th derivative $g^{(k)}$ of a function $g$ is defined by the recursive rule

$$g^{(0)}(x) = g(x) \qquad\qquad g^{(k+1)}(x) = \frac{d}{dx} g^{(k)}(x) \qquad k \geqslant 0.$$

More generally, given any set $X$, and any function $f : X \to X$, we define a sequence $(a_k)$ of elements of $X$ as follows.

$$a_1 \in X \quad \text{given;} \qquad\qquad a_{k+1} = f(a_k), \quad k \geqslant 1. \qquad\qquad (8.5)$$

The first term of the sequence (the initial condition) is given. Because $f$ is a function, if $a_k$ is well-defined for some $k \geqslant 1$, then so is $a_{k+1}$. The induction principle ensures that the whole sequence is well-defined. Since each element of $X$ gives rise to a sequence, we get lots of sequences with just one function!

The function $f$ may depend on several preceding terms

$$a_k = f(a_{k-1}, a_{k-2}, \ldots, a_{k-d})$$

in which case the integer $d$ is called the **order** of the sequence. In a recursive sequence of order $d$, the first $d$ terms of the sequence do not have the required number of preceding terms, so these terms must be supplied explicitly. In other words, there are $d$ initial conditions.

The **Fibonacci sequence** is a second-order sequence

$$a_0 = 1; \qquad a_1 = 1; \qquad a_k = f(a_{k-1}, a_{k-2}) = a_{k-1} + a_{k-2}, \quad k \geq 2.$$

We find

$$(1, 1, 2, 3, 5, 8, 13, 21, 34, 55, 89, 144, \ldots).$$

Because the term $a_k = a_{k+2} - a_{k+1}$ is a well-defined function of the following two terms, we can extend the Fibonacci sequence backwards, to obtain a doubly-infinite sequence

$$(\ldots, -21, 13, -8, 5, -3, 2, -1, 1, 0, 1, 1, 2, 3, 5, 8, 13, 21, \ldots).$$

A recursive sequence is said to be **invertible** if the function $f$ has an inverse at all points of the sequence. Thus the Fibonacci sequence is invertible.

The behaviour of a recursive sequence may be extremely complicated. Let $n$ be a natural number. We define a first-order recursive sequence

$$x_0 = n \qquad x_{t+1} = f(x_t), \quad t \geqslant 0 \tag{8.6}$$

by the simple rule

$$f : \mathbb{N} \to \mathbb{N} \qquad x \mapsto \begin{cases} x/2 & \text{if } x \text{ even} \\ 3x+1 & \text{if } x \text{ odd.} \end{cases}$$

The initial condition $x_0 = 1$ leads to a **periodic sequence**

$$(1,4,2,1,4,2,1,\ldots)$$

consisting of indefinite repetitions of the pattern $(1,4,2)$. If instead we start with $x_0 = 7$, we find

$$(7,22,11,34,17,52,26,13,40,20,10,5,16,8,4,2,1,4,2,1,\ldots)$$

so, after an irregular initial excursion, the sequence reaches the same periodic pattern as the previous sequence. The so-called '$3x+1$' **conjecture** states that

**Conjecture.** *For any n, the sequence (8.6) contains 1.*

This conjecture states that all these sequences are **eventually periodic**, and reach the same periodic pattern $(4,2,1)$.

Any open conjecture leads to functions whose existence cannot be decided. For example, we may define the function $\tau : \mathbb{N} \to \mathbb{N}$, where $\tau(n)$ is the smallest non-negative integer $t$ for which the sequence with initial condition $x_0 = n$ has $x_t = 1$. We compute a few values of $\tau(n)$.

| $n$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| $\tau(n)$ | 0 | 1 | 6 | 2 | 5 | 8 | 16 | 3 | 19 | 6 |

The behaviour of the function $\tau$ seems quite unpredictable; is this function well-defined? Millions of values of $n$ have been tested on computers, and all resulting sequences were found to have this property. However, nobody has been able to prove that this must happen for every $n$, and so, at present, we cannot claim that the function $\tau$ is well-defined (see exercise 7).

# Exercises

**Exercise 8.1.** Using the notation of section 8.1, write Dirichlet's box-principle and its negation symbolically.

**Exercise 8.2.** With reference to Dirichlet's box principle, determine the number of all possible arrangements of $m$ objects in $n$ boxes.

**Exercise 8.3.** Use Dirichlet's box principle to prove the following statements.

1. *Show that given any five points in the unit square, there are two points whose distance apart is at most $1/\sqrt{2}$.*
   [ *Divide the square into four suitable regions.*]

2. *Show that in any finite group of people, there are two people with the same number of friends.*

3. *Let $X$ be a finite set, and let $f : X \to X$ be a function. Show that every recursive sequence of elements of $X$ is eventually periodic.*

**Exercise 8.4.** Let $f$ be as in (8.3). Construct a proper subset of $\mathbb{N}$ that contains $f(\mathbb{N})$. (The smaller the set, the better.)

**Exercise 8.5.** Let $f$ be a positive and increasing real function. For every positive integer $n$, let $p_n$ be the smallest prime number $p$ in the range $n \leqslant p \leqslant n + f(n)$. Given $f$, the sequence $(p_n)$ may not be well-defined. Experiment with various functions $f$ to get an idea of what's needed to get a well-defined sequence. (The slower $f$ grows, the better —and the harder!)

**Exercise 8.6.** Consider the recursive sequence (8.5), where $f(x)$ is a quadratic polynomial.

1. Compute $a_k$ as a function of $a_1$ and $k$, for small values of $k$.

2. Write an essay on the behaviour of such a function, as $k$ increases.

**Exercise 8.7.** Consider the function $\tau$ associated to the '$3x+1$ conjecture' (section 8.3). Determine an infinite subset of $\mathbb{N}$ over which this function is well-defined. (The larger the set, the better.)

**Exercise 8.8.** Use the twin prime conjecture in order to define a function whose existence cannot be decided at present. Do the same with Goldbach conjecture.

# Chapter 9

# Bad mistakes

In constructing a mathematical argument it's easy to make mistakes. In this chapter we identify some common mistakes found in definitions, statements of theorems, and their proofs. Awareness of these problems should help us avoiding them.

There are, of courses, plenty of uninteresting mistakes: a word or symbol used in two senses, a deduction containing unnecessary steps, etc. Here we deal with more substantial problems, which render an argument invalid.

Mistaking examples for a proof.

Defining a non-existing object, or an ambiguous one.

Assuming what we are supposed to prove.

Applying a function to a point outside its domain.

Using a picture and failing to realise the geometrical assumptions involved.

## 9.1 Examples vs. proofs

Arguably, the most basic mistake in a mathematical argument is to believe that verifying the validity of a statement in specific examples constitutes a form of proof. Our study of Euler's polynomial in section 6.7 shows how misleading examples can be.

This state of affairs is peculiar to mathematics, and indeed to modern mathematics. Ancient Greek, Medieval Arab, and Renaissance western mathematicians, instead of giving general arguments, tended to offer examples to be copied. In other

scientific disciplines, e.g., biology, providing a 'proof' of a statement amounts to verifying its validity experimentally in a large number of cases.

Faulty arguments of this kind are often easy to spot, as in the following example.

**Theorem.**   *For all primes $p$, the integer $2^p - 2$ is divisible by $p$.*

BAD PROOF.

$$2^2 - 2 = 2 \cdot 1, \quad 2^3 - 2 = 3 \cdot 2, \quad 2^5 - 2 = 5 \cdot 6, \quad 2^7 - 2 = 7 \cdot 18, \quad \text{etc.} \quad \square$$

The theorem has been proved only for $p = 2, 3, 5, 7$.

Our next example is similar, but not so clear-cut (Suppes, *Introduction to Logic* p. 138f)

**Theorem.**   *For all $x$, $y$ and $z$, if $x + y < x + z$ then $y < z$.*

BAD PROOF.  Suppose $x + y < x + z$. Take $x = 0$. Then

$$y = 0 + y < 0 + z = z.$$

$\square$

The mistake here is that we took $x$ to be 0, which is a special value of $x$. This assumption is wholly unjustified, since all three quantities $x, y, z$ are controlled by an existential quantifier, and hence we are not allowed to impose any condition on them. By adding the assumption that $x = 0$, we have in fact proved the

WEAKER THEOREM. *For all $y$ and $z$, if $0 + y < 0 + z$, then $y < z$.*

This is not what we claimed to prove.


## 9.2   Bad definitions

An incorrect definition may imply that the object being defined is not the one we had in mind, or, worse, that it does not exist at all!

BAD DEFINITION. *Let $A$ be a finite set, and let $\mathbf{P}(A)$ be its power set. We define the function*

$$f : \mathbf{P}(A) \to \mathbb{Q} \qquad X \mapsto \frac{\#X}{\#A}$$

*which gives the relative size of a subset $X$ of $A$.*

This function is undefined if $A$ is empty! (For some, this observation exemplifies the kind of annoying pedantry mathematicians are notorious for.) This minor slip is fixed by writing '*Let $A$ be a finite non-empty set*'.

BAD DEFINITION. Let $X, Y, Z$ be sets and let

$$A := X \smallsetminus Y \smallsetminus Z.$$

The set $A$ exists, but it's not uniquely defined, because the set-difference operator is *not associative* (see section 2.1). In general, we have

$$(X \smallsetminus Y) \smallsetminus Z \neq X \smallsetminus (Y \smallsetminus Z)$$

so in the definition of $A$ we must clarify which of the two cases we have in mind.

Our next example shows an incorrect function definition that can be put right in many different ways.

BAD DEFINITION. *Let $f : \mathbb{Q} \to \mathbb{Q}$ be defined as follows*

$$f(x) = x^2 y + 1.$$

The problem here is rather obvious: the definition of $f(x)$ depends on another argument besides $x$, which is not specified. This mistake is fatal: as defined, $f$ is not a function. There are many legitimate interpretations of what this formula could mean, each resulting in a different function.

1. We re-define the domain, supplying the missing information via a second argument.

$$f : \mathbb{Q}^2 \to \mathbb{Q} \qquad f(x, y) = x^2 y + 1.$$

In this definition, $y$ too is rational. Clearly, there are many other possibilities.

2. The missing argument is regarded as a **parameter** (see section 5.1).

$$f_\lambda : \mathbb{Q} \to \mathbb{Q} \qquad f_\lambda(x) = \lambda x^2 + 1 \qquad \lambda \in \mathbb{Q}.$$

We have highlighted the change in status of the variable $y$ by switching to the Greek alphabet, and turning it into a subscript. For every value of $\lambda$, we have a different function.

3. The missing argument is regarded as an **indeterminate.** In this case $f(x)$ is a polynomial in $y$, and we must re-define the co-domain of $f$ accordingly.

$$f : \mathbb{Q} \to \mathbb{Q}[y] \qquad f(r) = r^2 y + 1.$$

The symbol $\mathbb{Q}[y]$ denotes the set of all polynomials in the indeterminate $y$, with coefficients in $\mathbb{Q}$ (see section 2.4.1). We have replaced $x$ with a symbol that is

normally not used for indeterminates, and which reminds us that this argument is rational.

In the next example, the function being defined requires more conditions than those given.

BAD DEFINITION. *Let $A, B, C$ be sets, and let $f : A \rightarrow C$ and $g : B \rightarrow C$ be functions. We define*

$$h : A \cap B \rightarrow C \qquad x \mapsto f(x) + g(x).$$

If $A$ and $B$ are disjoint, then the domain of the function $h$ is empty. Moreover, we have assumed implicitly that the elements of the co-domain can be added together. However, $C$ may be a set where addition is not defined at all (e.g., if $f$ and $g$ are predicates, then $C = \{T, F\}$, and adding boolean constants makes no sense), or $C$ may not be closed under addition, for instance, $C$ could be an interval. We give a more restrictive definition, where these problems are taken care of.

DEFINITION. *Let $A, B$ be sets with non-empty intersection, let $C$ be a set closed under addition, and let $f : A \rightarrow C$ and $g : B \rightarrow C$ be functions. We define*

$$h : A \cap B \rightarrow C \qquad x \mapsto f(x) + g(x).$$

The next definition hides a subtle problem.

BAD DEFINITION. *Let $x$ be a number in the closed unit interval, and let us consider its binary expansion*

$$x = \sum_{k=1}^{\infty} \frac{c_k}{2^k} \qquad c_k \in \{0, 1\}. \tag{9.1}$$

*We define the function $f$ that associates to $x$ its binary digits sequence*

$$f : [0, 1] \rightarrow \{0, 1\}^{\mathbb{N}} \qquad x \mapsto (c_1, c_2, \ldots).$$

What's wrong with this definition? The co-domain of $f$ is the set of all binary sequences, for which we have used the notation developed in section 5.2. By comparing the terms of the identity

$$\frac{1}{2^m} = \sum_{k=1}^{\infty} \frac{1}{2^{k+m}} \qquad m \in \mathbb{Z}$$

with equation (9.1), we see that the rational numbers $1/2^m$ have two distinct binary representations. For instance, letting $m = 1$, we have two distinct representations of $1/2$, namely

$$(1, 0, 0, 0, \ldots) \qquad \text{and} \qquad (0, 1, 1, 1, \ldots).$$

So our function is not uniquely defined at these rationals. We give an alternative definition, where the ambiguity is resolved by choosing consistently one of the two digit sequences.

DEFINITION. *Let $x$ be a number in the half-open unit interval $[0, 1)$. We consider the binary expansion of $x$*

$$x = \sum_{k=1}^{\infty} \frac{c_k}{2^k} \qquad c_k \in \{0, 1\}.$$

*Without loss of generality, we assume that, if the sequence $(c_k)$ is eventually constant, then such a constant is zero. We define the function $f$ that associates to $x$ its binary digit sequence*

$$f : [0, 1) \to \{0, 1\}^{\mathbb{N}} \qquad x \mapsto (c_1, c_2, \ldots).$$

In our next definition, we use a hidden assumption. Under such circumstances, using the definition adds extra information (we say that the definition is *creative*).

BAD DEFINITION. *In the set of congruence classes modulo an integer $m$, we define fractions by*

$$\frac{a + m\mathbb{Z}}{b + m\mathbb{Z}} = \text{the class } c + m\mathbb{Z} \text{ such that } cb \equiv a \, (\text{mod } m)$$

*where $b$ is not divisible by $m$.*

Take $m = 6$. We have $2 \times 3 = 0$, so $0/3 = 2$. But also $0 \times (1/3) = 0$, so

$$0 + 6\mathbb{Z} = \frac{0 + 6\mathbb{Z}}{3 + 6\mathbb{Z}} = 2 + 6\mathbb{Z}$$

but $0 \not\equiv 2 \, (\text{mod } 6)$. What went wrong? We have assumed that the implication

$$\text{if } a \neq 0 \text{ and } ab = ac, \text{ then } b = c,$$

which is valid for the integers, is also valid for congruence classes.

## 9.3   Bad implications

Inappropriate handling of implications is a common source of mistakes, and in many cases the problem originates from confusing an implication with its converse. For instance, from $a^2 = b^2$, deducing $a = b$, or from '$P \Rightarrow Q$' and $Q$, deducing $P$. Examples like these are sometimes called *non sequiturs* (Latin for 'doesn't follow').

An instance of this mistake is often found in 'proofs' of implications. This is a classic example.

> *Prove that* $\sqrt{2} + \sqrt{6} < \sqrt{15}$

BAD PROOF.

$$\sqrt{2} + \sqrt{6} < \sqrt{15} \;\Rightarrow\; (\sqrt{2} + \sqrt{6})^2 < 15$$
$$\Rightarrow\; 8 + 2\sqrt{12} < 15$$
$$\Rightarrow\; 2\sqrt{12} < 7$$
$$\Rightarrow\; 48 < 49$$

The last statement is true, which completes the proof.   $\square$

What's wrong with this proof? We must prove $P$, where $P = (\sqrt{2} + \sqrt{6} < \sqrt{15})$. Instead, we have assumed $P$, and correctly deduced a true statement. So we have proved

> $P \Rightarrow$ TRUE.

What do we know from this? Nothing, because both TRUE $\Rightarrow$ TRUE and FALSE $\Rightarrow$ TRUE are true statements, as one can see by checking the truth table for the operator $\Rightarrow$, given in (3.9). So the original statement $P$ could be true or false. For instance, had we started from the false statement $\sqrt{2} + \sqrt{6} < -\sqrt{15}$, we would have reached exactly the same conclusion.

There are two methods for fixing problems of this kind.

*First method: retracing the steps.*

We regard the chain of deductions displayed above as 'rough work'; then we start from the end and attempt to prove the chain of converse implications.

PROOF.

$$48 < 49 \;\Rightarrow\; \sqrt{48} < \sqrt{49}$$

$$\Rightarrow\ 2\sqrt{12} < 7$$
$$\Rightarrow\ 8 + 2\sqrt{12} < 15$$
$$\Rightarrow\ (\sqrt{2} + \sqrt{6})^2 < 15$$
$$\Rightarrow\ \sqrt{2} + \sqrt{6} < \sqrt{15}$$

where the first and the last implications are justified by the fact that we have taken the positive square root of each side. We have proved the proposition TRUE $\Rightarrow P$, from which we deduce that $P$ is TRUE. $\square$

Clearly, we could not have come up with such a proof without doing the 'rough work'.

*Second method: contradiction.*

PROOF. Let us assume that the statement to be proved is false

$$\neg P\ =\ (\sqrt{2} + \sqrt{6} \geqslant \sqrt{15})$$
$$\Rightarrow\ (\sqrt{2} + \sqrt{6})^2 \geqslant 15$$
$$\Rightarrow\ 8 + 2\sqrt{12} \geqslant 15$$
$$\Rightarrow\ 2\sqrt{12} \geqslant 7$$
$$\Rightarrow\ 48 \geqslant 49.$$

We have proved that $\neg P \Rightarrow$ FALSE, and hence $\neg P$ is necessarily FALSE, that is, $\neg(\neg P) = P$ is TRUE. $\square$

The advantage of this method is that it doesn't require any 'rough work'.

## 9.4   Bad use of functions

Common mistakes derive from applying a function to arguments outside its domain. These include dividing by zero (0 is not in the domain of $x^{-1}$), taking real square roots of negative numbers, forming $\tan(\pi/2)$, $\ln(0)$, etc.

Let $f, g$ be differentiable real function and let $h(x) = f(g(x))$. The chain rule for differentiating the composition of two functions states that

$$\frac{dh}{dx} = \frac{df}{dg}\frac{dg}{dx}.$$

Let us consider a proof of this statement.

BAD PROOF.    Let $\varepsilon$ be a real number. We compute

$$\frac{h(x+\varepsilon)-h(x)}{\varepsilon} = \frac{f(g(x+\varepsilon))-f(g(x))}{g(x+\varepsilon)-g(x)}\frac{g(x+\varepsilon)-g(x)}{\varepsilon}.$$

Letting $\varepsilon$ tend to zero gives the desired result.    $\square$

This argument has an appealing simplicity, but is also quite sloppy. What if $g$ is non-injective, so that $g(x+\varepsilon)=g(x)$ for some $\varepsilon \neq 0$? For instance, $g$ could be a constant function.

Our next problem originates from mishandling a non-invertible function, with rather amazing consequences!

BAD THEOREM.    *Every negative real number is positive.*

BAD PROOF.    Let $x$ be a negative real number. Then $x = -1 \cdot |x|$. Now

$$1 = 2 \cdot \frac{1}{2}$$

and hence

$$-1 = (-1)^1 = (-1)^{2 \cdot \frac{1}{2}} = [(-1)^2]^{\frac{1}{2}} = 1^{\frac{1}{2}} = 1.$$

Therefore $x = -1 \cdot |x| = 1 \cdot |x| = |x| > 0$, as required.    $\square$

What's wrong with this proof? The function $x \mapsto x^2$ is not one-to-one, and hence is not invertible. To invert it, we must restrict its domain. If we restrict the domain to the values $x \geq 0$, and take the positive sign for the square root, then $x \mapsto x^2$ and $x \mapsto \sqrt{x}$ are the inverse of each other: $x^{2 \cdot \frac{1}{2}} = \sqrt{x^2} = x$. However, if the argument $x$ is negative, to ensure that squaring and square root are still inverse of each other, one must take the negative sign of the square root.

Obviously, the mistake occurs in the chain of equalities that lead to the statement $-1 = 1$. The faulty step is the last one

$$1^{\frac{1}{2}} = 1,$$

where the original negative argument $x = -1$ is paired with the positive sign of the square root. Had we chosen the correct negative sign, we would have recovered the value $-1$ that we started with. Exercise 4 deals with a similar problem.

## 9.5 Bad pictures

When using a picture in an argument, one must be careful that in the picture there are no hidden geometric assumptions.

BAD THEOREM. *Every triangle is isosceles.*

BAD PROOF. Take any triangle *ABC*. We show that $|AB| = |AC|$.
For contradiction, suppose $|AB| \neq |AC|$. Draw the line *n* bisecting the angle *BAC*, and the perpendicular bisector *m* of the side *BC*. Since $|AB| \neq |AC|$, the lines *n* and *m* are not parallel, and hence they meet, say at a point *P*.

Drop perpendiculars *PF* and *PE* from *P* to *AB* and *AC*. Now the angles *FAP*, *EAP* are equal and the triangles *FAP*, *EAP* have *AP* in common, so $|FP| = |EP|$ and $|AF| = |AE|$. Also $|BD| = |CD|$, so $|BP| = |CP|$. Since the angles *BFP* and *CEP* are equal, it follows that $|BF| = |CE|$. So

$$|AB| = |AF| + |BF| = |AE| + |CE| = |AC|.$$

So the sides *AB* and *AC* do have equal length, contradicting our assumption. □

The problem with this proof is that we assumed *P* lies inside the triangle. In fact it never does, and only one of the two perpendiculars *PF*, *PE* from *P* has an end in the triangle.

Several proofs in Euclid's Elements, even of true theorems, make unjustified assumptions of this kind. Some people draw the conclusion that one should never use pictures in mathematical arguments, because they might contain hidden assumptions. This is silly.

## Exercises

**Exercise 9.1.** The following definition has several flaws. (*a*) Explain what they are; (*b*) write a correct, clearer definition.

Let $a = (b_1, b_2, \ldots)$ be a given sequence of elements of B. We define the function

$$f : \mathbb{Z} \to B \qquad f(m) = \sum_{k=1}^{m} b_k^{-1}.$$

**Exercise 9.2.** The following theorem and proof have several faults. (*a*) Explain what they are; (*b*) write an appropriate revision.

BAD THEOREM.    *For all numbers a and b,*

$$\frac{a^2 + b^2}{|ab|} > 2.$$

BAD PROOF.    For

$$\frac{a^2 + b^2}{|ab|} > 2 \;\Rightarrow\; a^2 + b^2 > 2ab$$
$$\Rightarrow\; a^2 - 2ab + b^2 > 0$$
$$\Rightarrow\; (a - b)^2 > 0.$$

The last inequality is trivially true, which proves it.    □

**Exercise 9.3.** The following statements are false, and in their inductive 'proof' there is a flaw. Explain clearly and concisely what the flaw is. (You shouldn't need more than one sentence.)

(a) BAD THEOREM.    *For any natural number n, the following holds*

$$2 + 4 + \cdots + 2n = (n - 1)(n + 2).$$

BAD PROOF.    Assume that $2 + 4 + \cdots + 2k = (k - 1)(k + 2)$ for some $k \in \mathbb{N}$. Then

$$
\begin{aligned}
2 + 4 + \cdots + 2(k + 1) &= (2 + 4 + \cdots + 2k) + 2(k + 1) \\
&= (k - 1)(k + 2) + 2(k + 1) \\
&\qquad \text{(by the induction hypothesis)} \\
&= k^2 + k - 2 + 2k + 2 \\
&= k(k + 3) \\
&= [(k + 1) - 1)][(k + 1) + 2].
\end{aligned}
$$

which is the given statement for $n = k + 1$. It follows that the statement is true for all $n \in \mathbb{N}$.    □

(b) BAD THEOREM. (Pólya[1])    *For all $n \in \mathbb{N}$, in any group of n horses, all horses have the same colour.*

---

[1] George Pólya, Hungarian mathematician (1887–1985).

BAD PROOF.    For any natural number $n$, we consider the statement

$$\mathscr{P}(n) = \text{'in any set of } n \text{ horses all horses have the same colour'}$$

The proposition $\mathscr{P}(1)$ is clearly true. Assume now that $\mathscr{P}(n)$ is true for some $n \geq 1$, and consider an arbitrary collection of $n+1$ horses



By the induction hypothesis, the first $n$ horses have the same colour, and so do the last $n$ horses.



But then all $n+1$ horses have the same colour as the $n-1$ horses common to the two sets.



This completes the proof.    □

**Exercise 9.4.** Identify the flaw in the proof of the following statement, and explain it in great detail.

BAD THEOREM.    *Trigonometry does not exist. Specifically, for all $\theta \in \mathbb{R}$, we have $\cos(\theta) = 1$ and $\sin(\theta) = 0$.*

BAD PROOF.    Let $\theta$ be a real number. We find

$$e^{i\theta} \;=\; e^{i\theta \frac{2\pi}{2\pi}} = e^{2\pi i \frac{\theta}{2\pi}} = \left(e^{2\pi i}\right)^{\frac{\theta}{2\pi}} = 1^{\frac{\theta}{2\pi}} = 1.$$

From the formulae

$$\cos(\theta) = \frac{e^{i\theta} + e^{-i\theta}}{2} \qquad\qquad \sin(\theta) = \frac{e^{i\theta} - e^{-i\theta}}{2i}$$

the result follows.    □

[*The flaw is subtle: it's related to an incorrect inversion of a function. In section section 9.4 we dealt with a similar problem.* ]

# Bibliography

[1] Harley Flanders, *Amer. Math. Monthly* (1971) 1–10.

[2] Leonard Gillman, *Writing mathematics well,* The Mathematical Association of America, Washington, D.C. (1987) [ISBN: 0-88385-443-0].

[3] R. P. Boas, Can we make mathematics intelligible? *Amer. Math. Monthly* (1971) 727–731.

[4] H. Cohn, *Advanced number theory*, Dover, New York (1980).

[5] N. J. Higham *Handbook of Writing for the Mathematical Sciences* SIAM, Philadelphia (1998).

[6] K. Houston, *How to think like a mathematician,* Cambridge University Press, Cambridge (2009).

[7] Steven G. Krantz, *A primer of Mathematical Writing,* American Mathematical Society (1997).

[8] Donald E. Knuth, Tracy L. Larrabee, and Paul M. Roberts *Mathematical writing,* The Mathematical Association of America, Washington, D.C. (1989) [ISBN 0-88385-063-X].

[9] William Strunk and E. B. White, *The elements of style,* fourth edition, Longman, New York (1999) [ISBN: 0-205-30902-X].

[10] Gérard Tenenbaum and Michel Mendès France, *The prime numbers and their distribution*, Student Mathematical Library, Volume 6, American Mathematical Society, Providence, Rhode Island (2000).

[11] Lynne Truss, *Eats, shoots and leaves: the zero-tolerance approach to punctuation,* Gotham (Penguin) (2004).

[12] Mary-Claire van Leunen, *A Handbook for Scholars,* Oxford University Press (1992).

[13] *Writing matters,* the Royal Literary Fund report on student writing in higher education, Edited by Stevie Davis, David Swinburne, and Gweno Williams, The Royal Literary Fund, London (2006). (Electronic copy available from `http://www.rlf.org.uk/fellowshipscheme/research.cfm`.)

# Index