

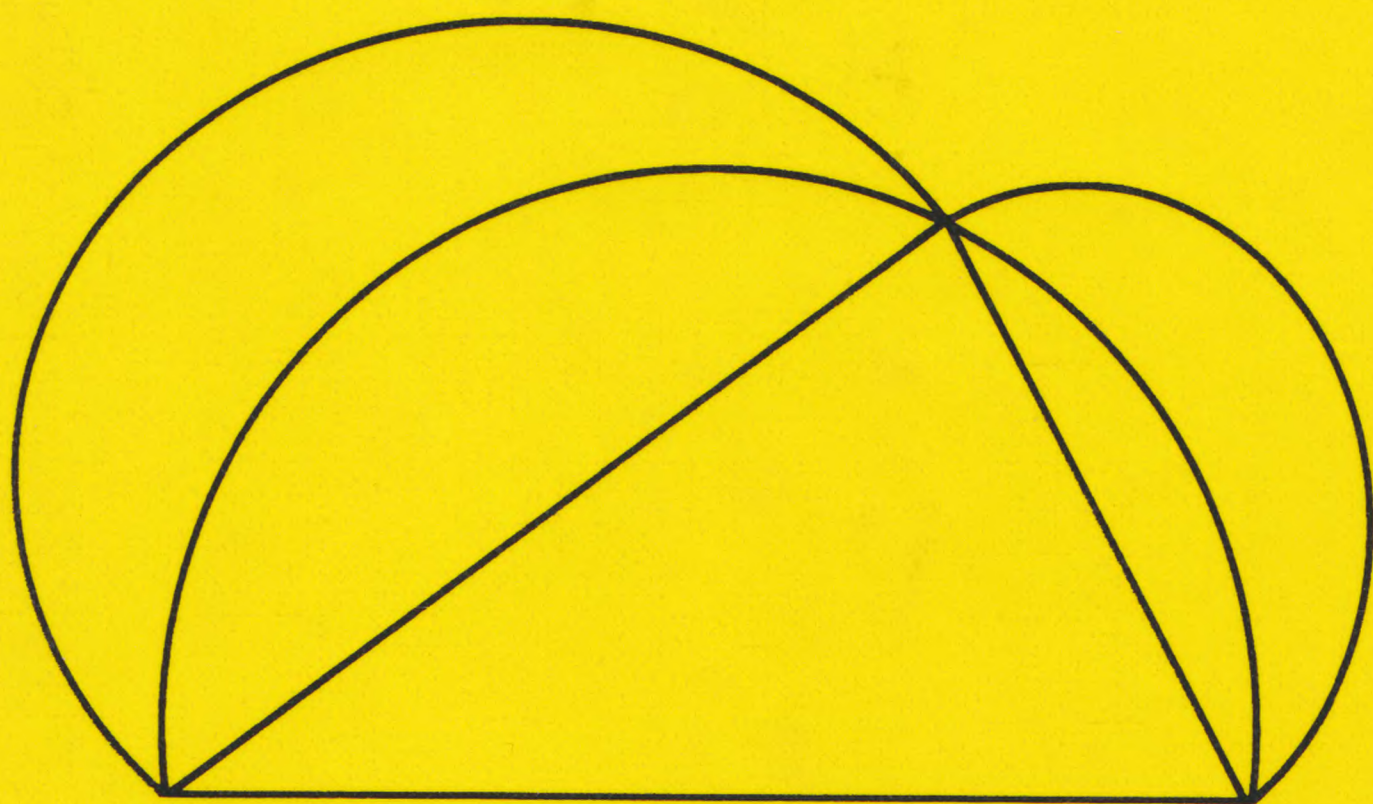
Undergraduate Texts in Mathematics

Readings in Mathematics

W.S. Anglin

J. Lambek

The Heritage of Thales



Springer

Undergraduate Texts in Mathematics

Readings in Mathematics

Editors

S. Axler

F.W. Gehring

P.R. Halmos

Springer

New York

Berlin

Heidelberg

Barcelona

Budapest

Hong Kong

London

Milan

Paris

Tokyo

Graduate Texts in Mathematics

Readings in Mathematics

Ebbinghaus/Hermes/Hirzebruch/Koecher/Mainzer/Neukirch/Prestel/Remmert: *Numbers*

Fulton/Harris: *Representation Theory: A First Course*

Remmert: *Theory of Complex Functions*

Undergraduate Texts in Mathematics

Readings in Mathematics

Anglin: *Mathematics: A Concise History and Philosophy*

Anglin/Lambek: *The Heritage of Thales*

Bressoud: *Second Year Calculus*

Hämmerlin/Hoffmann: *Numerical Mathematics*

Isaac: *The Pleasures of Probability*

Samuel: *Projective Geometry*

W.S. Anglin
J. Lambek

The Heritage of Thales

With 23 Illustrations



Springer

W. S. Anglin
J. Lambek
Department of Mathematics
and Statistics
McGill University
Montreal, Quebec
Canada H3A 2K6

Editorial Board

S. Axler
Dept. of Mathematics
Michigan State University
East Lansing, MI 48824
USA

F.W. Gehring
Dept. of Mathematics
University of Michigan
Ann Arbor, MI 48109
USA

P.R. Halmos
Dept. of Mathematics
Santa Clara University
Santa Clara, CA 95053
USA

Mathematics Subject Classification (1991): 01-01, 01A05

Library of Congress Cataloging-in-Publication Data
Anglin, W.S.

The heritage of Thales / W.S. Anglin and J. Lambek.
p. cm. — (Undergraduate texts in mathematics. Readings in
mathematics)

Includes bibliographical references and index.

ISBN 0-387-94544-X (hc : alk. paper)

1. Mathematics—History. 2. Mathematics—Philosophy. I. Lambek,
Joachim. II. Title. III. Series.

QA21.A535 1995

510'.9—dc20

95-19695

Printed on acid-free paper.

© 1995 Springer-Verlag New York, Inc.

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher (Springer-Verlag New York, Inc., 175 Fifth Avenue, New York, NY 10010, USA), except for brief excerpts in connection with reviews or scholarly analysis. Use in connection with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed is forbidden.

The use of general descriptive names, trade names, trademarks, etc., in this publication, even if the former are not especially identified, is not to be taken as a sign that such names, as understood by the Trade Marks and Merchandise Marks Act, may accordingly be used freely by anyone.

Production managed by Natalie Johnson; manufacturing supervised by Jeffrey Taub.

Photocomposed copy prepared from the authors' LaTeX file.

Printed and bound by R.R. Donnelley & Sons, Harrisonburg, VA.

Printed in the United States of America.

9 8 7 6 5 4 3 2 1

ISBN 0-387-94544-X Springer-Verlag New York Berlin Heidelberg

Preface

This is intended as a textbook on the history, philosophy and foundations of mathematics, primarily for students specializing in mathematics, but we also wish to welcome interested students from the sciences, humanities and education. We have attempted to give approximately equal treatment to the three subjects: history, philosophy and mathematics.

History

We must emphasize that this is not a scholarly account of the history of mathematics, but rather an attempt to teach some good mathematics in a historical context. Since neither of the authors is a professional historian, we have made liberal use of secondary sources. We have tried to give references for cited facts and opinions. However, considering that this text developed by repeated revisions from lecture notes of two courses given by one of us over a 25 year period, some attributions may have been lost. We could not resist retelling some amusing anecdotes, even when we suspect that they have no proven historical basis. As to the mathematicians listed in our account, we admit to being colour and gender blind; we have not attempted a balanced distribution of the mathematicians listed to meet today's standards of political correctness.

Philosophy

Both authors having wide philosophical interests, this text contains perhaps more philosophical asides than other books on the history of mathematics. For example, we discuss the relevance to mathematics of the pre-Socratic philosophers and of Plato, Aristotle, Leibniz and Russell. We also have

presented some original insights. However, on some points our opinions diverge; so, in a spirit of compromise, we have agreed to excise some of our more extreme views. Some of these divergent opinions have been expressed in Anglin [1994] and Lambek [1994].

Mathematics

One of the challenges one faces in offering a course on the history and philosophy of mathematics is to persuade one's colleagues that the course contains some genuine mathematics. For this reason, we have included some mathematical topics, usually not treated in standard courses, for example, the renaissance method for solving cubic equations and an elementary proof of the impossibility of trisecting arbitrary angles by ruler and compass constructions. We have taken the liberty of presenting many mathematical ideas in modern garb, with the hindsight inspired by more recent developments, since a presentation faithful to the original sources, while catering to the serious scholar, would bore most students to tears.

In Part I we deal essentially with the history of mathematics up to about 1800. This is because thereafter mathematics tends to become more specialized and too advanced for the students we have in mind. However, we make occasional excursions into more modern mathematics, partly to relieve the tedium associated with a strictly chronological development and partly to present modern answers to some ancient questions, whenever this can be done without overly taxing the students' ability.

In Part II we deal with some selected topics from the nineteenth and twentieth centuries. In that period, mathematics became rather specialized and made spectacular progress in different directions, but we confine attention to questions in the foundations and philosophy of mathematics.

The more universal aspects of mathematics are sketched briefly in the last five sections. We introduce the language of category theory, which attempts a kind of unification of different branches of mathematics, albeit at a very basic and abstract level.

Acknowledgements

The authors wish to acknowledge partial support from the Social Sciences and Humanities Research Council of Canada, from the Natural Sciences and Humanities Research Council of Canada and from the Quebec Department of Education.

We wish to express our sincerest thanks to Matthew Egan for his undaunted dedication in typing and editing and to Mira Bhargava, Henri Darmon and Ramona Behravan for their conscientious reading and criticism of the manuscript.

W. S. Anglin and J. Lambek

Contents

Preface	v
0 Introduction	1
PART I: History and Philosophy of Mathematics	5
1 Egyptian Mathematics	7
2 Scales of Notation	11
3 Prime Numbers	15
4 Sumerian-Babylonian Mathematics	21
5 More about Mesopotamian Mathematics	25
6 The Dawn of Greek Mathematics	29
7 Pythagoras and His School	33
8 Perfect Numbers	37
9 Regular Polyhedra	41
10 The Crisis of Incommensurables	47
11 From Heraclitus to Democritus	53

12 Mathematics in Athens	59
13 Plato and Aristotle on Mathematics	67
14 Constructions with Ruler and Compass	71
15 The Impossibility of Solving the Classical Problems	79
16 Euclid	83
17 Non-Euclidean Geometry and Hilbert's Axioms	89
18 Alexandria from 300 BC to 200 BC	93
19 Archimedes	97
20 Alexandria from 200 BC to 500 AD	103
21 Mathematics in China and India	111
22 Mathematics in Islamic Countries	117
23 New Beginnings in Europe	121
24 Mathematics in the Renaissance	125
25 The Cubic and Quartic Equations	133
26 Renaissance Mathematics Continued	139
27 The Seventeenth Century in France	145
28 The Seventeenth Century Continued	153
29 Leibniz	159
30 The Eighteenth Century	163
31 The Law of Quadratic Reciprocity	169
 PART II: Foundations of Mathematics	 173
1 The Number System	175
2 Natural Numbers (Peano's Approach)	179
3 The Integers	183

4	The Rationals	187
5	The Real Numbers	191
6	Complex Numbers	195
7	The Fundamental Theorem of Algebra	199
8	Quaternions	203
9	Quaternions Applied to Number Theory	207
10	Quaternions Applied to Physics	211
11	Quaternions in Quantum Mechanics	215
12	Cardinal Numbers	219
13	Cardinal Arithmetic	223
14	Continued Fractions	227
15	The Fundamental Theorem of Arithmetic	231
16	Linear Diophantine Equations	233
17	Quadratic Surds	237
18	Pythagorean Triangles and Fermat's Last Theorem	241
19	What Is a Calculation?	245
20	Recursive and Recursively Enumerable Sets	251
21	Hilbert's Tenth Problem	255
22	Lambda Calculus	259
23	Logic from Aristotle to Russell	265
24	Intuitionistic Propositional Calculus	271
25	How to Interpret Intuitionistic Logic	277
26	Intuitionistic Predicate Calculus	281
27	Intuitionistic Type Theory	285

28 Gödel's Theorems	289
29 Proof of Gödel's Incompleteness Theorem	291
30 More about Gödel's Theorems	293
31 Concrete Categories	295
32 Graphs and Categories	297
33 Functors	299
34 Natural Transformations	303
35 A Natural Transformation between Vector Spaces	307
References	311
Index	321

0

Introduction

Remarks on prehistory

Long before written records were kept, people were concerned with the seasons, important in agriculture, and the sky, which permitted them to read off the passage of time. Everyone knows that the *year* is the time it takes the sun to complete its orbit about the earth. (Copernicus notwithstanding, mathematical readers will see nothing wrong with placing the origin of the coordinate system at the center of the earth.) Also, a *month* is supposed to be the time it takes the moon to go around the earth; at least, this was the case before the lengths of the months were laid down by law. But what about the *week*? Theological explanations aside, it is the smallest period, longer than a day, that can be easily observed by looking at the sky: the time it takes the moon to pass from one phase to another, from new moon to half moon, from half moon to full moon, etc.

The days of the week are named after the sun, the moon and the five planets visible to the naked eye: Mars (French *mardi*), Mercury (French *mercredi*), Jupiter (French *jeudi*), Venus (French *vendredi*) and Saturn (English *Saturday*). The English *Tuesday*, *Wednesday*, *Thursday* and *Friday* are named after the Teutonic deities which supposedly correspond to the Roman gods after whom the planets were named.

In Hindu astronomy there are nine planetary deities, the *graha*. In addition to the seven associated with the days of the week, there are two others, *rahu* and *kebu*, alleged to be associated with the so-called ‘nodes’. These are the points where the orbits of the sun and the moon, when traced out

on the firmament, intersect. (See Freed and Freed [1980].) The importance of the nodes is that an eclipse of the sun or the moon can only occur when both sun and moon fall on the nodes, to within 10° ; according to an ancient rule of thumb, this happens once in about 18.6 years.

At Stonehenge in England there is an imposing prehistoric monument, dating from about 2,500 BC. The huge standing stones of the monument were presumably used to sight the points on the horizon where the sun and the moon, and perhaps Venus, rise and set at certain dates (Hawkins [1965]). They are surrounded by a circle of 56 holes in the ground, and Fred Hoyle [1977] has proposed the ingenious hypothesis that these were used as a calendar and to calculate the dates of possible eclipses.

According to Hoyle, the idea was to move a *sun marker* two holes in 13 days, a *moon marker* two holes each day, and two *nodal stones* three holes per year. The sun marker would thus complete an orbit in 364 days; the discrepancy could be fixed by appropriate adjustments at midsummer and midwinter. The moon marker would complete an orbit in 28 days, that is, four weeks. Of course, this should really be 29.5 days, so adjustments might have to be made each full moon and each new moon. The nodal stones would take $56/3$ years to perform a complete orbit. On those occasions when both sun marker and moon marker were about to catch up with the nodal stones, the presiding priest could risk predicting an eclipse.

Foreword on history

Even so-called 'primitive' societies may be engaged in some fairly sophisticated mathematical activities, for example, the calculations involved in kinship descriptions. (How many students can tell on the spot what exactly is a *second cousin three times removed*?) For interested readers, we recommend two recent books: *Africa Counts* by Claudia Zaslowsky and *Ethnomathematics* by Marcia Ascher.

Mathematics, as usually conceived, begins with the development of agriculture in the river valleys of Egypt, Iraq, India and China. If we pay more attention to the Near East than to the Far East, this is because the former has provided us with more accessible records and because modern mathematics can be traced back directly to it. We possess written records concerning the state of mathematics in Egypt and Mesopotamia (Iraq) from as early as about 2000 BC. Around 500 BC, mathematical knowledge spread to the Greek world. This included not only modern Greece, but also the coast of Asia Minor (modern Turkey) and Magna Grecia (southern Italy and Sicily). About 300 BC, the center of mathematics moved from Athens to Alexandria in Egypt, where it was to remain for the next 800 years; for it was in Alexandria that all the books were kept.

Around 500 AD, mediterranean civilization finally came to a stop, per-

haps because of the repeated impact of epidemic diseases. About 800 AD, mathematics in the Alexandrian tradition resurfaced in India, which had a long mathematical tradition of its own. The Arabs, aided by translations of Greek texts, developed and transmitted mathematical knowledge from India back to the mediterranean area and ultimately to Europe. During the so-called 'renaissance', mathematics flourished in Italy and, aided by the Chinese invention of printing, spread to Western and Central Europe. Of course, today mathematics is being pursued in all the industrial countries of the world.

Introduction to the number system

The historical and pedagogical development of the number system goes somewhat like this:

$$\mathbf{N}^+ \rightarrow \mathbf{Q}^+ \rightarrow \mathbf{R}^+ \rightarrow \mathbf{R} \rightarrow \mathbf{C} \rightarrow \mathbf{H} .$$

Here \mathbf{N}^+ is the set of positive integers, the *numbers* used for counting, known to all societies. \mathbf{Q}^+ is the set of positive rationals, namely, *quotients* of positive integers, surely known to all agricultural civilizations. At one time, they were believed to exhaust all the numbers, until the Pythagoreans discovered that the diagonal of a square was not a rational multiple of its side. We use \mathbf{R}^+ to denote the positive *reals*; these were certainly used effectively by Thales, though the Greeks originally tended to regard them as ratios of geometric quantities. A formal treatment, anticipating the nineteenth century definition by Dedekind, was first given by Eudoxus in Athens. The transition from \mathbf{R}^+ to \mathbf{R} , the set of all reals, positive, zero and negative, took place in India and may be ascribed to Brahmagupta. The set \mathbf{C} of *complex* numbers was first considered by Cardano to describe the intermediate steps in solving a cubic equation with real coefficients and three real solutions. The set \mathbf{H} of *quaternions* is named after their inventor William Hamilton, who may have been preceded by Olinde Rodrigues and perhaps even by Carl Friedrich Gauss.

Most of the advances in the development of the number system may have been motivated by the desire to solve equations. Thus, the equations $2x = 3$, $x^2 = 2$, $x + 1 = 0$ and $x^2 + 1 = 0$ led to the successive introduction of \mathbf{Q}^+ , \mathbf{R}^+ , \mathbf{R} and \mathbf{C} , respectively. However, all polynomial equations with complex coefficients do have complex solutions, so the introduction of quaternions requires a different justification. They were motivated by the desire to pass from the plane, describable by complex numbers, to three or four dimensions.

Part I
Topics in the History and
Philosophy of Mathematics

1

Egyptian Mathematics

The Greeks believed that mathematics originated in Egypt. As to the reason for this, opinion was divided. Aristotle thought that mathematics was developed by priests, ‘because the priestly class was allowed leisure’ (*Metaphysics* 981b 23-24). Herodotus believed that the annual flooding of the Nile necessitated surveying to redetermine field boundaries, and thus led to the invention of geometry. In fact, Democritus referred to Egyptian mathematicians as ‘rope stretchers’. It may be of interest to note that the Egyptians themselves believed that mathematics had been given to them by the god Thoth. Our only original sources of information on the mathematics of ancient Egypt are the Moscow Mathematical Papyrus and the Rhind Mathematical Papyrus.

The Moscow Papyrus dates from 1850 BC, about the time the Bible dates the life of the patriarch Abraham. In 1893 it was acquired by V. S. Golenishchev and brought to Moscow (Gillings [1972], p. 246). Problem 14 of this papyrus is by far the most interesting. It is the computation of a *truncated pyramid*, a square pyramid with a similar pyramid cut off its top. If a side of the base has length a and a side of the top has length b , then the volume of the truncated pyramid of vertical height h is

$$V = \frac{h}{3}(a^2 + ab + b^2).$$

This is exactly the formula used by the Egyptians. Note that, if $b = 0$, we get the formula for the volume of the complete pyramid.

The Rhind Mathematical Papyrus seems to be based on an earlier work. It was written by one Ahmose in 1650 BC, about the time when, accord-

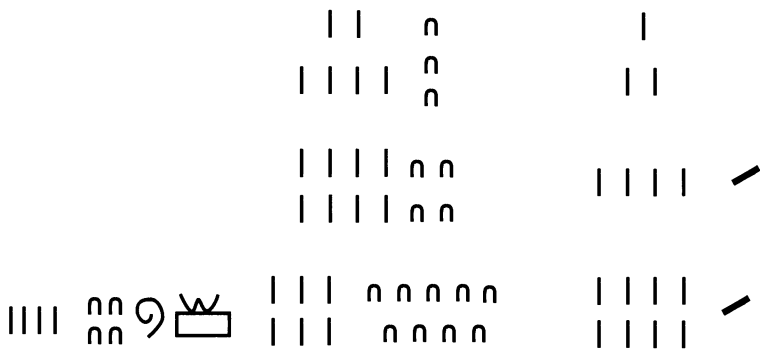


FIGURE 1.1. Rhind Papyrus

ing to the Bible, Joseph was governor of Egypt. Alexander Henry Rhind acquired it in Luxor, Egypt in 1858; the British Museum bought it from his estate in 1865. Complete photographs of the papyrus can be found in *The Rhind Mathematical Papyrus* edited by G. Robins and C. Shute.

The Rhind Papyrus opens by promising the reader ‘a thorough study of all things, insight into all that exists, knowledge of all obscure secrets’. It is a bit of a letdown to find that it is, in fact, a sequence of solved problems in elementary mathematics, a sort of Schaum’s outline for aspiring scribes. These scribes had to calculate how many bricks were needed to build a ramp of a certain size, how many loaves of bread were required to feed the labourers, and so on.

Problem 32 of the papyrus is an exercise in multiplication written as in Figure 1.1.

Transcribing this into modern notation, we have

12	1	
24	2	
48	4	/
96	8	/

144 = the sum of the checked entries.

Clearly, this is a calculation to show that $12 \times 12 = 144$, using the fact that $12 = 4 + 8$.

By doubling and adding, the Egyptians were able to multiply any two natural numbers – without having to memorize multiplication tables! Sometimes they used a different, yet equivalent, method, as illustrated by the following multiplication of 70 by 13:

70	13	/
140	6	
280	3	/
560	1	/

910 = sum of checked entries.

We let the reader figure out why this works. The method of repeated doubling can also be used for division. In the following example, we divide

184 by 17 (stopping at 136, as the next double exceeds 184):

17	1	
34	2	/
68	4	
136	8	/

The Egyptians would first check off the last row and subtract 136 from 184, obtaining 48. They would then check off the row containing 34, the highest multiple of 17 less than 48. Since $48 - 34 = 14$ is less than 17, they would now add up all the entries in the second column with check marks beside them: $2 + 8 = 10$. This gives the answer: the quotient is 10 and the remainder is 14.

In carrying out these divisions, the Egyptians sometimes interspersed doubling with multiplication by 10 (their language expressed numbers in the base 10, just as ours does). For example, Problem 69 in the Rhind Mathematical Papyrus is to calculate the number of 'ro' of flour in each loaf, if 1120 ro of flour is made into 80 loaves. In other words, we are asked to divide 1120 by 80:

80	1	
800	10	/
160	2	
320	4	/

sum of checked numbers = 14.

The Egyptians also knew how to extract square roots and how to solve linear equations. They used the hieroglyph h much as we use the letter x for the unknown. They used the formula $(\frac{4}{3})^4 r^2$ for the area of a circle (which implies 3.16 as an approximation to π) and they did some interesting work with arithmetic progressions. For example, Problem 64 of the Rhind Papyrus is to find an arithmetic progression with 10 terms, with sum 10, and with common difference $1/8$.

In using fractions, the Egyptians were hampered by a curious tradition. They insisted on expressing all fractions (except $2/3$) as the sum of distinct *unit* fractions of the form $1/n$, n being a positive integer. Thus $2/9$ would be written as $1/6 + 1/18$ and $19/8$ as $2 + 1/4 + 1/8$. Even $2/3$ is sometimes written as $1/2 + 1/6$.

For us it is easy to divide $5/13$ by 12, but for the Egyptians this was a substantial problem. To help with such problems, they had a table listing unit fraction decompositions for fractions of the form $2/n$, with n an odd positive integer. This table (found in the Rhind Papyrus) gives $2/13$ as $1/8 + 1/52 + 1/104$. Since $5 = (2 \cdot 2) + 1$, Ahmose would write

$$\begin{aligned} 5/13 &= 2(1/8 + 1/52 + 1/104) + 1/13 \\ &= 1/4 + 1/26 + 1/52 + 1/13. \end{aligned}$$

From this he would obtain

$$(5/13)/12 = 1/48 + 1/312 + 1/624 + 1/156.$$

Actually, any fraction of the form $2/(2m+1)$ can be expressed as a sum of the unit fractions $1/(m+1)$ and $1/(m+1)(2m+1)$. Note that the Egyptians always followed this recipe; for example, Ahmose wrote $2/45 = 1/30 + 1/90$.

Recently, Paul Erdős proposed the following problem: show that, if n is an odd integer greater than 4, then $4/n$ can be written as a sum of three distinct unit fractions. The problem has not yet been solved. (See Mordell, p. 287.)

Exercises

1. Derive the formula for the volume of a truncated pyramid from that of a pyramid.
2. Explain why the above method for multiplying 70×13 works.
3. Find two ways of writing $1/4$ as the sum of two distinct unit fractions.
4. If m is a positive integer, show that $4/(4m+3)$ can be written as the sum of three distinct unit fractions.

2

Scales of Notation

The ancient Egyptian language belongs to the Hamito-Semitic group of languages. Like the Indo-European group, it contains a system of counting by tens, undoubtedly arising from the habit of counting using one's fingers. The notation used for writing numbers is also clearly based on the scale of ten. For some reason, standard French departs from decimal nomenclature; it expresses 97 as 4 times 20 plus 17. This seems odd, since it was the French who introduced the decimal system for weights and measures.

Some African languages express numbers in the scale of five. One may express natural numbers in any scale b , where b is an integer greater than 1, since every natural number is uniquely expressible in the form

$$a = a_0 + a_1b + a_2b^2 + a_3b^3 + \cdots + a_nb^n,$$

where $0 \leq a_i < b$ (for $i = 0, 1, \dots, n$). We write this more briefly as

$$a = (a_na_{n-1} \cdots a_2a_1a_0)_b.$$

If there is no doubt which scale is in use, the subscript b may be dropped.

The Egyptians had a number system based on the scale of ten, but, as we saw above, they often worked with scale two: to multiply by 12, Ahmose expressed 12 as $4 + 8$, that is, $2^2 + 2^3$, or $12 = (1100)_2$. The Egyptians also took $b = 7$ in some of their calculations (Gillings, p. 227), since there are seven *palms* in a *cubit*. They had no symbol for zero; instead they used special symbols for different powers of ten.

The binary scale (with $b = 2$) shows up in the Chinese *Book of Changes* (1200 BC), a system of divination in which each six place binary number

represents some concept. The digit 1 was associated with the male 'yang', and the digit 0 with the female 'yin'. The number $34 = (100010)_2$ was supposed to represent 'progress and success'. The binary scale also shows up in the Hindu classification of meters in verse, about 800 BC. Finally, it is of course used in the modern computer. The digit 1 is represented by a current, and the digit 0 by the absence of a current. Number scales are often found in recreational mathematics, as in the following three problems.

Six Weight Problem: A balance is a weighing apparatus with a central pivot, a beam, two scales and a set of counter-weights that are placed in one of the scales. Suppose we have some flour and we want to be able to put it into bags weighing anywhere from one to sixty-three kilograms. How can this be done using just six counter-weights?

Answer: Weights of 1, 2, 4, 8, 16 and 32 kilograms will allow you to weigh any integral load from 1 to $32 + 16 + 8 + 4 + 2 + 1 = 63$ kilograms.

Four Weight Problem: This time, suppose we are allowed to put weights on either scale. How can we weigh bags under 42 kilograms using only four weights?

Answer: Choose weights of 1, 3, 9 and 27 kilograms, since any integer a can be written uniquely in the form

$$a = a_0 + a_1 3 + a_2 3^2 + \cdots + a_n 3^n,$$

where each a_i is one of -1 , 0 or 1 .

The Game of Nim: This so-called Chinese game is played by two opponents, who take turns removing matches from several piles according to the following rules:

1. A player must remove at least one match in a turn.
2. A player may remove any number of matches from a single pile in a turn.

The player who removes the last match wins. Find a strategy for winning this game.

Answer: Express the number of matches in each pile in the scale of two and write these binary numbers one below the other. If, when it is your turn, you can arrange it so that each column adds up to an even number, then you can do the same in every subsequent turn and you will win the game.

As an example, suppose there are three piles containing 7, 5 and 3 matches. It is your turn. In binary notation, the piles contain the following number of matches:

$$\begin{array}{ccc} 1 & 1 & 1 \\ 1 & 0 & 1 \\ & 1 & 1 \end{array}$$

To make the number of 1's in each column even, you take a match from the first pile, leading to

$$\begin{array}{ccc} 1 & 1 & 0 \\ 1 & 0 & 1 \\ & 1 & 1 \end{array}$$

Your opponent takes 2 matches, say, from the third pile, leaving

$$\begin{array}{ccc} 1 & 1 & 0 \\ 1 & 0 & 1 \\ & 1 & \end{array}$$

You take two matches from the first pile, yielding

$$\begin{array}{ccc} 1 & 0 & 0 \\ 1 & 0 & 1 \\ & 1 & \end{array}$$

Your opponent then removes all the matches from the first pile, say, resulting in

$$\begin{array}{ccc} 1 & 0 & 1 \\ & 1 & \end{array}$$

You now take 4 matches from the first pile, leaving just one match in each pile. Your opponent has to take one of them, and you win by taking the last match.

What is going on in the Six Weight Problem? It is easily seen that the following three statements are equivalent:

1. Every integral load less than 64 kg can be weighed uniquely with 6 weights: 1, 2, 4, 8, 16 and 32 kg.
2. Every natural number less than 64 can be expressed uniquely as the sum of distinct powers of 2.
3. Every natural number less than 64 can be written uniquely in the scale of 2 with up to 6 digits, each 0 or 1.

A direct proof is quite easy, but a proof in the spirit of the 18th century is more interesting. In preparation for this proof, let us look at the following multiplication:

$$(1+x^2)(1+x^3)(1+x^5)(1+x^7) = 1+x^2+x^3+2x^5+2x^7+x^8+x^9+2x^{10}+\dots$$

Why are the coefficients of x^5 , x^4 and x^3 equal to 2, 0 and 1, respectively? Because $5 = 2 + 3$ can be written as the sum of some of 2, 3, 5 and 7 in two ways (the first sum consists of one term only), 4 cannot be written as the sum of some of 2, 3, 5 and 7 at all, and 3 can only be written as the sum of these numbers in one way, the sum having one term. In general, the coefficient of x^n will be the number of ways in which n can be written as the sum of distinct members of the set $\{2, 3, 5, 7\}$.

Now let us replace this set of numbers by the set $\{1, 2, 4, 8, 16, 32\}$ and consider

$$(1+x)(1+x^2)(1+x^4)(1+x^8)(1+x^{16})(1+x^{32}) = \sum_{i=0}^{\infty} f(n)x^n.$$

Then $f(n)$ is the number of ways in which n can be written as the sum of distinct powers of 2, up to 32. Clearly, $f(n) = 0$ when $n \geq 64$. What if $n < 64$? The left-hand side can also be written

$$\begin{aligned} \frac{1-x^2}{1-x} \cdot \frac{1-x^4}{1-x^2} \cdot \frac{1-x^8}{1-x^4} \cdot \frac{1-x^{16}}{1-x^8} \cdot \frac{1-x^{32}}{1-x^{16}} \cdot \frac{1-x^{64}}{1-x^{32}} &= \frac{1-x^{64}}{1-x} \\ &= 1 + x + x^2 + \cdots + x^{63} = \sum_{n=0}^{63} x^n. \end{aligned}$$

Hence $f(n) = 1$ if $n < 64$.

Suppose, instead of stopping with x^{32} , we form the *infinite* product $\prod_{n=0}^{\infty} (1+x^{2^n})$. Then we can show similarly that *every* natural number can be written in the scale of 2.

Exercises

1. Write out a scale 7 multiplication table.
2. Show how to convert a scale 10 numeral to a scale 7 numeral.
3. Give a proof, in the spirit of the 18th century, that every natural number can be written uniquely in the scale of 3. (Hint: form the infinite product $(1+x+x^2)(1+x^3+x^6)(1+x^9+x^{18})\cdots$ and evaluate it in two different ways.)
4. Likewise, show that any integer can be written uniquely as $\sum_{k=0}^n a_k 3^k$, where $a_k = -1, 0$ or 1 .

3

Prime Numbers

It would be impossible to write a history of mathematics without mentioning prime numbers, and it would be improper to give an account of prime numbers without going into the history of mathematics. Prime numbers enter into almost every branch of mathematics; they are as fundamental as they are ubiquitous. Their history can be used as a framework for a history of mathematics generally. In this chapter, we take a brief look at the fascinating subject of primes.

The Egyptians might have written

$$\frac{4}{5} = \frac{1}{2} + \frac{1}{4} + \frac{1}{20}.$$

From this, it follows that

$$\frac{4}{10} = \frac{1}{4} + \frac{1}{8} + \frac{1}{40}$$

and that

$$\frac{4}{15} = \frac{1}{6} + \frac{1}{12} + \frac{1}{60}.$$

The moral to be drawn from this is that, to express a/b as a sum of unit fractions, it suffices to consider the case when b cannot be factored into smaller numbers. An integer greater than 1 which cannot be factored into numbers, all of which are smaller than the original integer, is called *prime*. The first few primes are

$$2, 3, 5, 7, 11, 13, 17, 19, 23, 29, 31, 37, \dots$$

Note that a positive integer is *prime* if and only if it has exactly two positive integer divisors.

Early on, people noticed that a pile of small stones can sometimes be arranged in a rectangle and sometimes it cannot. Thus, although we do not have any record of this, the Egyptians probably knew the difference between composite and prime numbers. Indeed, it is not impossible that some Egyptian scribe may have noticed that, if every proper fraction of the form $4/p$, with p prime, and greater than 3, can be expressed as a sum of three distinct unit fractions, then every proper fraction of the form $4/n$, with n any positive integer greater than 4, can be so expressed. (See the problem of Erdős, mentioned in Chapter 1.)

It was the Greeks who first *proved* that the number of primes is infinite. A proof is found in Euclid's *Elements* (300 BC).

Euclid's Lemma (Book VII Proposition 31):

Every integer $n > 1$ is divisible by some prime number.

Proof: Among the divisors of n which are greater than 1, let p be the smallest. Then p has no divisors other than 1 and p — any other divisor of p would be a divisor of n as well — and hence p is prime.

Euclid's Theorem (Book IX Proposition 20):

Given any finite list of primes $p_1, p_2 \dots p_k$, there is a prime not on this list.

Proof: Consider the number $n = p_1 p_2 \dots p_k + 1$. Clearly, n is not divisible by any of the primes on the list; for, upon dividing n by p_i , we get remainder 1. From the lemma we know that n does have a prime factor q (possibly n itself). Hence there is a prime, namely, q , which is not on the list. QED.

In Proposition 14 of Book IX, Euclid proved that, if n is a square-free positive integer (that is, one with no square factor other than 1), then n has a factorization into primes which is unique (if you list the prime factors in order of increasing size). However, it was not until 1801 that the unique factorization was formally proved for *any* positive integer n . This was done by Carl Friedrich Gauss (1777–1855) in his *Disquisitiones Arithmeticae*. Although mathematicians used the unique factorization theorem long before 1801, and although almost any one of them could have found a proof for it, Gauss was the first person actually to sit down and do so. Perhaps the other mathematicians considered the theorem too obvious to be worth proving. One way to prove that every positive integer greater than 1 has a unique factorization into primes is as follows.

Proof of the Unique Factorization Theorem:

Let n be the smallest positive integer, if there is one, which has 2 (or more) factorizations into primes:

$$n = pqr \dots = p'q'r' \dots$$

We assume the primes are written in nondecreasing order. By minimality of n , $p \neq p'$ (or we could cancel off the p 's and get a smaller number with

two factorizations). Without loss of generality, we may suppose that $p' < p$. Hence

$$p' < p \leq q \leq r \leq \dots \quad (*)$$

Since n is not prime, $n \geq p^2$, and hence $n > pp'$. By minimality of n , $n - pp'$ has a unique factorization. Both p and p' are factors of $n - pp'$ (since $n - pp' = p(qr \cdots - p') = p'(q'r' \cdots - p)$) and hence, for some positive integer z , $n - pp' = pp'z$. This gives $qr \cdots - p' = p'z$, so that p' is a factor of $qr \cdots$. Since $qr \cdots < n$, $qr \cdots$ has a unique factorization into primes. Thus p' is one of q, r, \dots . But this contradicts $(*)$ above. For another proof, see Part II, Chapter 15.

Like Euclid, Eratosthenes of Cyrene (230 BC) worked at the University of Alexandria. He suggested a method for making a list of all prime numbers, which is called the 'sieve of Eratosthenes'. His method is as follows: write down all the positive integers greater than 1; cross out all multiples of 2 other than 2, cross out all multiples of 3 other than 3 which have not been crossed out yet, etc. In the end, the numbers not crossed out from a complete list of primes.

People often wonder whether there is a simple formula representing prime numbers. For example, $f(x) = x^2 - x + 41$ is prime for all integer values of x from 0 to 40. While this might convince a physicist that $f(x)$ is always prime, unfortunately $f(41) = 41^2$.

In 1743, Christian Goldbach observed that a polynomial

$$f(x) = a_0 + a_1x + a_2x^2 + \cdots + a_nx^n$$

with integer coefficients a_0, a_1, \dots, a_n cannot represent primes only, that is, the integers $f(0), f(1), f(2), \dots$ are not all prime.

Indeed, if $f(0) = p$, then $f(kp)$ is clearly a multiple of p for all integers k . But, as k tends to infinity, so does the absolute value of $f(kp)$. Hence, for some value of k , $f(kp)$ will be a proper multiple of p and therefore not prime.

It therefore came as a great surprise to the mathematical community when, in 1970, Yuri Matiyasevič formed a polynomial $f(x, y, z, \dots)$ with integer coefficients, but in several variables, such that, when positive integers are chosen for x, y, z, \dots , one gets all the prime numbers and only the prime numbers as positive values of the polynomial. We shall say more about this in Chapter 21 on Hilbert's Tenth Problem in Part II.

In 1830 (in *Théorie des Nombres* Vol. II, p. 65), A. M. Legendre noted that, if $\pi(x)$ is the number of primes less than or equal to x , then $\pi(x)$ is approximately equal to $x/(\log_e x - 1.08366)$, where $e = \lim_{n \rightarrow \infty} (1 + \frac{1}{n})^n$ is the base of the natural logarithm (Chapter 26). We shall write $\log x$ and assume the base to be e . He was not able to prove this. In 1896, two mathematicians working independently proved that

$$\lim_{x \rightarrow \infty} \frac{\pi(x)}{x/\log x} = 1.$$

These two mathematicians were the Frenchman Jacques Hadamard (1865–1963) and the Belgian Charles Jean de la Vallée Poussin (1866–1962). The result they proved is called the *Prime Number Theorem*. It implies that the n th prime is approximately equal to $n \log n$. For, if we let p_n be the n th prime, the equation implies that n is roughly equal to $p_n / \log p_n$, so that

$$p_n \approx n \log p_n \approx n \log(n \log n) \approx n \log n,$$

since $n \log \log p_n$ can be neglected in comparison with $n \log p_n \approx p_n$.

It was Goldbach who conjectured that every even number greater than 2 is a sum of two primes. This conjecture has not yet been proved or disproved. However, in 1937, the Russian mathematician I. M. Vinogradov made some progress towards proving Goldbach's Conjecture, by showing that every odd integer greater than, say, $10^{10^{10}}$ (or some similar bound) is a sum of three prime numbers. Some progress in the Goldbach conjecture was recently made by the Chinese mathematician Chen Jing-Run. He proved that every sufficiently large (say $> 10^{10^{10}}$) even number has the form $p + q$, where p is prime and q is either prime or the product of two primes. During the so-called 'cultural revolution' in the sixties this kind of mathematics was frowned upon in China for being far removed from any conceivable application to industry or agriculture. Because he stubbornly stuck to his esoteric research at the risk of neglecting his teaching, Chen Jing-Run was discriminated against during the reign of the so-called 'gang of four' and may have lost his academic position. After the overthrow of the gang of four, he was rehabilitated and even declared a 'hero of the revolution'.

At the moment (1995), one of the 'hot topics' in prime number theory is cryptography. In its simplest form, the idea is this: the cipher key is a product $n = pq$ of two large primes, typically having 50 to 80 digits each. Knowing n is enough to encode messages, but decryption requires knowledge of the factorization. The integer n is made public (hence the term 'public key') so that everyone can use the code to encipher messages. Security is maintained, because only the intended recipient knows the key, namely, the factorization pq , necessary to carry out the decryption. The basis for this scheme is that it takes a very long time to factor products of large primes and the war may well be over before the enemy succeeds in doing so. (Try to factor the relatively small product 1,315,685,447, and you will see that the enemy does not have an easy task.)

At the moment, much research is being done to find refinements of the above idea, refinements that are at once economical and secure for those who want to send secret messages. Much research is also being done to find ways of using computers to factor very large numbers, and thus break the codes based on the above idea.

Exercises

1. Find the 25 primes less than 100 and express 100 as the sum of two primes.
2. Prove that there exist 1,000 consecutive positive integers none of which is prime. (Hint: start with $1001! + 2$.)
3. Prove that there are infinitely many prime numbers of the form

$$4m - 1.$$

(Hint: consider $n = 4q_1q_2q_3 \cdots q_k - 1$, where the q_i are primes of the form $4m - 1$, and show that not all prime divisors of n can be of the form $4m + 1$.)

Sumerian-Babylonian Mathematics

The Sumerians were a people of unknown linguistic affinity, who lived in the southern part of Mesopotamia (Iraq), and whose civilization was absorbed by the Semitic Babylonians around 2000 BC. Babylonian culture reached its peak in about 575 BC, under Nebuchadnezzar, but most of the mathematical achievements we shall discuss in this chapter and in Chapter 5 are much older, going back as far as 2000 BC — about the time when the biblical patriarch Abraham was said to have been born in the Sumerian city of Ur.

As we shall see, Mesopotamian mathematics is quite impressive. However, we should remember that, like the ancient Egyptians, the Mesopotamians never gave what we would call ‘proofs’ for their results; the first people to do so were the Greeks.

In representing numbers up to (and including) 59, the Sumerians and Babylonians used a decimal system. For example, they wrote 35 as follows, where we have approximated the original cuneiform figures by ours:

$$\begin{array}{r} <<< & YYY \\ & YY \end{array}$$

On the other hand, 60 is again denoted by Y , and so is 60^2 , as well as 60^{-1} , 60^{-2} , etc. It is usually clear from the context which is meant. Here are some further examples:

$$<<< = 30, \text{ or } 30/60 = 1/2;$$

$$\begin{aligned}
 & < YY &= 12, \text{ or } 1/5; \\
 Y < < \begin{array}{c} Y \quad Y \\ Y \end{array} &= 84, \text{ or } 7/5.
 \end{aligned}$$

The Babylonian use of scale 60 was taken over into Greek astronomy around 150 BC by Hipparchus of Nicaea and it is still used today in measuring time and angles. To remove ambiguities in the above three examples, we would write

$$\begin{aligned}
 &30^\circ \text{ or } 30', \\
 &12^\circ \text{ or } 12', \\
 &1^\circ 24' \text{ or } 1^\circ 24''.
 \end{aligned}$$

The scale 60, or *sexagesimal system*, was also employed for weights of silver: 60 shekels = 1 mina; 60 minas = 1 talent. The prophet Ezekiel, living in Babylon, wrote in 573 BC:

The Lord Yahweh says this: ... Twenty shekels, twenty-five shekels and fifteen shekels are to make one mina (*Ezekiel* 45:9–12).

The later Babylonians even introduced a symbol for zero:

$$Y \leq \begin{array}{c} Y \quad Y \\ Y \end{array} = 60^2 + 4 = 3604.$$

Ptolemy (150 AD) replaced this symbol by a small circle, probably from the Greek word ‘ouden’, meaning ‘nothing’.

In order to divide, the Babylonians made use of the fact that $a/b = a \cdot b^{-1}$. To this end, they constructed tables of inverses, like the one given in Table 4.1 (taken from Neugebauer [1969]). Note that the scribe did not list the inverses of any integers having a prime factor other than 2, 3 or 5. It seems he was afraid of repeating sexagesimals!

The Babylonians also had tables of squares, cubes, square roots, cube roots, and even roots of the equations

$$x^2(x \pm 1) = a.$$

Their method for extracting square roots is sometimes called *Heron's method* after Heron of Alexandria (60 AD), who included it in his *Metrica*. Let a_1 be a rational number between \sqrt{a} and $\sqrt{a} + 1$, where a is a positive non-square integer; let $a_{n+1} = (a_n + a/a_n)/2$; then $a_n \rightarrow \sqrt{a}$ as $n \rightarrow \infty$. Indeed, if $e = a_1 - \sqrt{a}$, we have $0 < e < 1$ and

$$0 < a_{n+1} - \sqrt{a} < 2\sqrt{a}(e/2\sqrt{a})^{2^n}$$

(see Exercise 4). As $n \rightarrow \infty$, this tends to 0.

b	1/b	b	1/b
2	30'	16	3'45''
3	20'	18	3'20''
4	15'	20	3'
5	12'	24	2'30''
6	10'	25	2'24''
8	7'30'	27	2'13''20''
9	6'40'	30	2'
10	6'	32	1'52''30''
12	5'	36	1'40''
15	4'	40	1'30''

TABLE 4.1. Mesopotamian table of inverses (scale 60)

For example, if $a = 2$, $a_1 = 3/2$, then $a_2 = 17/12$ and $a_3 = 577/408$. In sexagesimal notation, $577/408 = 1^\circ 24' 51'' 10''' 35'''' \dots$. The fourth approximation $a_4 = 665857/470832$, which is $1^\circ 24' 51'' 10''' 7'''' \dots$ in sexagesimal notation. The difference between a_4 and $\sqrt{2}$ is less than

$$2\sqrt{2} \left(\frac{3/2 - \sqrt{2}}{2\sqrt{2}} \right)^{2^4} < 10^{-23}.$$

The Babylonian tablet YBC7289, dating from about 1600 BC, gives $\sqrt{2}$ as $1^\circ 24' 51'' 10'''$.

Exercises

1. Write 5000 in the Babylonian manner. (You may use our degrees, minutes and seconds.)
2. Let a/b be a proper, reduced fraction (with a and b positive integers). Let $e_1 = 60a/b$ and $e_{n+1} = 60(e_n - [e_n])$ – where $[e_n]$ is the greatest integer less than or equal to e_n . Prove that the Babylonian sexagesimal expansion for a/b is

$$([e_1][e_2][e_3] \dots)_{60}.$$

3. Express $1/7$ as a repeating sexagesimal.
4. Prove by mathematical induction that

$$0 < a_{n+1} - \sqrt{a} < 2\sqrt{a}(e/2\sqrt{a})^{2^n}.$$

5. Use the Babylonian method to find $\sqrt{3}$ to within 60^{-10} .

6. Let a/b be a proper, reduced fraction (with a and b positive integers). Prove that a/b has an infinitely repeating sexagesimal expansion if and only if b has a prime factor which does not divide 60.

5

More about Mesopotamian Mathematics

In *Science Awakening I*, B. L. van der Waerden quotes the beginning of ‘AO8862’, a Babylonian clay tablet going back to about the same time as the Rhind Papyrus:

Length, width. I have multiplied the length and the width, thus obtaining the area. Then I added to the area, the excess of the length over the width: 183 was the result. Moreover, I have added the length and the width: 27. Required length, width and area.

27 and 183, the sums; 15 the length; 180 the area; 12 the width;

One follows this method:

$$27 + 183 = 210, 2 + 27 = 29.$$

Take one half of 29 (this gives $14\frac{1}{2}$),

$$14\frac{1}{2} \times 14\frac{1}{2} = 210\frac{1}{4},$$

$$210\frac{1}{4} - 210 = \frac{1}{4}.$$

The square root of $\frac{1}{4}$ is $\frac{1}{2}$.

$$14\frac{1}{2} + \frac{1}{2} = 15, \text{ the length;}$$

$$14\frac{1}{2} - \frac{1}{2} = 14, \text{ the width.}$$

Subtract 2, which has been added to 27, from 14, the width. 12 is the actual width. I have multiplied the length 15 by the width 12.

$$15 \times 12 = 180, \text{ the area;}$$

$$15 - 12 = 3;$$

$$180 + 3 = 183.$$

What is going on here? In modern notation, we would write x and y for length and width, respectively. The problem is to find a solution for the simultaneous equations

$$xy + (x - y) = 183 \text{ and } x + y = 27.$$

The answer is given as $x = 15$ and $y = 12$. The scribe's method is this: consider

$$xy + x - y + x + y = x(y + 2) = 210.$$

Putting $y' = y + 2$, we have $xy' = 210$. On the other hand, adding the factors of 210, we get

$$x + y' = x + y + 2 = 29;$$

$$\text{hence } \frac{1}{2}(x + y') = \frac{1}{2}(29) = 14\frac{1}{2};$$

$$\text{hence } \frac{x^2 + 2xy' + y'^2}{4} = (14\frac{1}{2})^2 = 210\frac{1}{4};$$

$$\text{hence } \frac{x^2 - 2xy' + y'^2}{4} = 210\frac{1}{4} - 210 = \frac{1}{4} \text{ (the so-called discriminant);}$$

$$\text{hence } \frac{x - y'}{2} = \frac{1}{2}.$$

Adding and subtracting $\frac{1}{2}(x + y')$ and $\frac{1}{2}(x - y')$, we get $x = 14\frac{1}{2} + \frac{1}{2} = 15$ and $y' = 14\frac{1}{2} - \frac{1}{2} = 14$. Note that 14 is not really the width; but $y = y' - 2 = 14 - 2 = 12$ is. The scribe then computes the area and checks his work. The scribe did not consider the possibility $x = 14, y + 2 = 15$, which gives the second solution $x = 14, y = 13$. He did not know how to take the negative square root of $\frac{1}{4}$.

The Babylonians could solve many kinds of equations, including: $ax = b$, $x^2 \pm ax = b$, $x^3 = a$, $x^2(x + 1) = a$. They could also solve simultaneous equations having the following forms:

$$x \pm y = a, \quad xy = b;$$

$$x \pm y = a, \quad x^2 + y^2 = b.$$

They even managed to solve the following pair of equations:

$$x^3\sqrt{x^2 + y^2} = 3,200,000; \quad xy = 1200. \quad (*)$$

As we saw just above, the Babylonians knew that

$$a^2 - b^2 = (a + b)(a - b).$$

They also knew that

$$(a + b)^2 = a^2 + 2ab + b^2.$$

Like the Egyptians, the Babylonians built pyramids, or *ziggurats*. If each story of a ziggurat consists of a square platform measuring $1 \times m \times m$, then the volume of a ziggurat with a base of length n is

$$(1 \times n \times n) + \cdots + (1 \times 2 \times 2) + (1 \times 1 \times 1) = 1^2 + 2^2 + 3^2 + \cdots + n^2.$$

The Babylonians knew that the formula for this sum is

$$n(n+1)(2n+1)/6,$$

a result also known to Pythagoras, but perhaps first proved by Archimedes.

According to the biblical story of the Tower of Babel, there was once an attempt to build a ziggurat 'with its top reaching heaven'. Perhaps the people behind this project thought that there was only a finite distance between heaven and earth, or perhaps they thought that they could calculate the sum of $1^2 + 2^2 + 3^2 + \cdots$, not realizing that the series diverges.

A remarkable fact about ancient Babylonian mathematics is that it included not just the so-called theorem of Pythagoras, but a theory of 'Pythagorean triangles'. (A *Pythagorean triangle* is a triple (x, y, z) of positive integers such that $x^2 + y^2 = z^2$, and thus x, y and z are sides of a right angled triangle.) From a clay tablet called 'Plimpton 322' (dating from 1900–1600 BC), we can deduce that the Babylonians used a result of which the following is a modern version:

Suppose u and v are *relatively prime positive integers*, that is, integers whose greatest common divisor is 1. Assume that not both are odd and that $u > v$. Then, if $a = 2uv, b = u^2 - v^2$ and $c = u^2 + v^2$, we have $\gcd(a, b, c) = 1$ and $a^2 + b^2 = c^2$.

Included in Plimpton 322 is the triangle (13500, 12709, 18541), which is generated by taking $u = 125$ and $v = 54$.

The converse of the above theorem is also true. That is, if a, b and c are relatively prime positive integers, with a even, such that $a^2 + b^2 = c^2$, then there are relatively prime positive integers u and v , not both odd, such that $a = 2uv, b = u^2 - v^2$ and $c = u^2 + v^2$. It is not impossible that the Babylonians knew this, but the earliest record we have of this result is in the solutions of Problems 8 and 9 of Book II of the *Arithmetica* of Diophantus (250 AD). Indeed, since Diophantus explained his ideas in terms of special cases, it is correct to say that the first explicit, rigorous proof of the converse of the Babylonian theorem was given only in 1738, by C. A. Koerber (Dickson [1971], Vol. II).

According to a tablet found in 1936 in Susa, an ancient city in what is now Iran, the Babylonians sometimes used the value $3\frac{1}{8}$ for π . At other times, they seem to have been satisfied with $\pi \approx 3$. It has been suggested that this Babylonian usage is behind 1 *Kings* 7:23–24:

He [Hiram of Tyre] made the basin of cast metal, ten cubits from rim to rim, circular in shape and five cubits high; a cord

thirty cubits long gave the measurement of its girth. Under its rim and completely encircling it were gourds; they went around the basin over a length of thirty cubits.

But perhaps the basin was hexagonal and not circular!

Exercises

1. Consider the simultaneous equations $xy + x - y = a$ and $x + y = b$, where a and b are given integers. What is a necessary and sufficient condition on a and b so that x and y will be integers?
2. Solve the simultaneous pair (*) (from a Susa tablet).
3. Prove by mathematical induction that

$$1^2 + 2^2 + 3^2 + \cdots + n^2 = n(n+1)(2n+1)/6.$$

4. Rabbi Nehemiah (150 AD) was unhappy with the idea that the Bible had used a very inaccurate value for π and he suggested that ‘the diameter of 10 cubits included the walls of the basin, while the circumference excluded them.’ Assuming that he was right, and assuming that the Bible used a perfectly accurate value for π , how wide was the wall (or rim) of the basin?
5. Prove that if a triangle has sides of lengths a, b and c , and if $a^2 + b^2 = c^2$, then the triangle is right angled.
6. Prove the Babylonian theorem for Pythagorean triangles.
7. Prove the converse of the Babylonian theorem for Pythagorean triangles.
8. In 1901, L. Kronecker gave the first proof that all positive integer solutions of $a^2 + b^2 = c^2$ are given without duplication by $a = 2uvk$, $b = (u^2 - v^2)k$, $c = (u^2 + v^2)k$, where u, v and k are positive integers such that $u > v$, u and v are not both odd, and u and v are relatively prime. Prove Kronecker’s theorem.
9. The 15 Pythagorean triangles in Plimpton 322 have angles which approximate the 15 whole number angles from 44° to 58° inclusive. Find a Pythagorean triangle, with relatively prime sides, one of whose angles is within $2/5^\circ$ of 47° . (Hint: see Anglin [1988].)

The Dawn of Greek Mathematics

The ancient Greek world was not confined to what we now call Greece, but extended to Ionia (western Turkey) in the east, southern Italy and Sicily in the west, and later to Alexandria (Egypt) in the south. Not surprisingly, Greek philosophy and mathematics began in Ionia, where the influence of the older civilizations of the east (e.g., Babylon) was greatest. Later, political events caused many Greeks to emigrate from Ionia to Italy, and this became the center of intellectual life for a while. After the war between a Greek coalition and the Persians ended in the defeat of the latter (490 BC), philosophy and mathematics flourished in Athens. Ultimately, after the founding of Alexandria (331 BC), it was there that most of the major scientific developments took place until about 500 AD.

The first Greek mathematician and philosopher is Thales of Miletus (600 BC). According to Proclus, Thales visited Egypt and brought back the knowledge of geometry from there. He may also have been influenced by Indian thought via Persia. He is said to have predicted the solar eclipse which occurred over the Near East in May of 585 BC. To do this, he may have made use of observations which the Babylonians had accumulated over many centuries.

Plato repeats a story about Thales being an absent-minded professor, who was so preoccupied with celestial matters that he did not observe what was in front of his feet and once fell into a well (*Theaetetus* 174a). According to other anecdotes, however, Thales could turn his mind to practical matters when necessary. He constructed an almanac and he used the theory of similar triangles to calculate the distance of ships from shore and

the height of pyramids. To impress his business minded fellow citizens, he once cornered the market in olive oil and, incidentally, made himself rich.

Thales's name is associated with a number of elementary theorems in geometry:

1. a circle is bisected by a diameter;
2. the base angles of an isosceles triangle are equal;
3. when two lines intersect, vertically opposite angles are equal;
4. the angle-angle-side congruence theorem;
5. the angle subtended by a diameter of a circle is a right angle (that is, if A, B, C are points on a circle and AC is a diameter, then $\angle ABC$ is a right angle).

Theorem 5 is called *Thales's theorem*. To prove this he also had to know the following:

6. the sum of the angles in a triangle is equal to two right angles (or, as we now say, in slavish imitation of the Babylonians, 180°).

All of these theorems must have been known empirically by the Egyptians and Babylonians. The reason they are associated with Thales is not that he discovered them, but that he was the first to prove them. This was the essential difference between pre-Greek and Greek mathematics: the Greeks established the logical connections among their results; they gave the first abstract proofs in mathematics.

As a philosopher, Thales is known for his statement that everything is made of water. How should we interpret this, and what is its relevance to mathematics?

As we look around us, we observe that there are two kinds of things: those that can be counted, such as pebbles and cows, and those that can only be measured, such as butter and water. This physical distinction between 'discrete' and 'continuous' is reflected on a linguistic level: it is perfectly correct to say 'one cow, two cows' but it sounds rather odd to say 'one butter, two butters', to put it mildly (however, see Exercise 5). We call the former sort of nouns *count nouns*, and the latter *mass nouns*. To some extent, this distinction is a convention. For example, in modern English, we can count peas but not rice, while a hundred years ago, 'pease' was not a plural but a mass noun. (A hundred years from now, 'rice' may be the plural of 'rouse'.)

A question which physicists are still working on is this: is the material universe ultimately countable — consisting of discrete, unconnected fragments — or is the material universe ultimately continuous — that is, should it be understood in terms of a connected continuum? If the first, how do we explain the unity of nature, how do we understand the continuity of change and motion? If the second, how do we explain the diversity of nature, how do we understand the individuality of distinct, single objects?

This issue was addressed by more than one Greek thinker. As we shall see, Pythagoras and Democritus took the view that reality is basically discrete. They then tried to understand apparently continuous entities in terms of discrete entities (e.g., lengths as ratios of whole numbers). Thales, on the other hand, took the view that ‘all is water’. In other words, the material universe is best understood in terms of a single substance, namely, water. (Here we are using the word ‘substance’ not in the Aristotelian sense of ‘individual entity’, but in the more common sense of ‘material having uniform properties’.) Thales had undoubtedly noticed that ice and steam are both forms of water, but we do not know why he picked water as the fundamental substance. (It has been suggested by Marxist historians that this was so because his city Miletus was a *maritime* power.) What is important is not that Thales overlooked the possibility of, say, there being 90 different substances, but that he raised a fascinating problem about the universe, which has not been resolved to this day.

Other Ionian philosophers agreed with Thales that there was a single substance, but not that it was water. Anaximenes of Miletus (550 BC) identified the primal substance as air. Heraclitus of Ephesus (500 BC) held that everything is made of fire. Anaximander was a follower and compatriot of Thales, who like Thales, took the view that the universe is best understood in terms of a single substance. Unlike Thales, he did not think this substance was water. He thought it was something he called the *Infinite*. The Infinite could take on the forms of earth, water, air, and fire. Today we might refer to solid, liquid, gas and energy, respectively.

Exercises

1. Let ABC be a triangle, and let d be a straight line through A parallel to BC . Assuming that the ‘alternate angles are equal’, prove that the sum of the angles of ABC equals two right angles.
2. Prove the Theorem of Thales, using Exercise 1 and the theorem that the base angles of an isosceles triangle are equal.
3. Prove the converse of Thales’s Theorem: if A , B and C are points on a circle and $\angle ABC$ is a right angle then AC is a diameter.
4. How would you measure the height of a pyramid (or tree, for that matter), using similar triangles?
5. Some nouns like ‘rice’ are definitely count nouns. We cannot say ‘two rices’. Other nouns are more problematic. ‘Whisky’ is normally a mass noun, but ‘two whiskies, please’ is perfectly good English (meaning two glasses of whisky). Some languages have many fewer count nouns than English. For example, in Indonesian it is incorrect to say ‘two

cows'; you have to say 'two tails of cow', as we might say 'two head of cattle'. (It is amusing to note that our mass noun 'cattle' is itself ultimately derived from the Latin word 'caput' meaning 'head'.) Write an essay on the distinction between count nouns and mass nouns and its relevance to mathematics.

Pythagoras and His School

Pythagoras (570–500 BC) was born in Samos, a Greek island off the coast of what is now Turkey. According to ancient sources (Iamblichus, Porphyry and Diogenes Laërtius), he traveled and studied in the Persian empire, which extended then from northern Greece to the Indus Valley and included ancient Mesopotamia. We know (Plimpton 322) that the Babylonians understood what is now called the ‘theorem of Pythagoras’, although the latter may have given the first proof. Pythagoras may have learned the theory of ‘Pythagorean triangles’ from the Babylonians.

According to the above mentioned sources, Pythagoras also studied under the Zoroastrian priests, the so-called ‘Magi’. However, judging from his belief in reincarnation and his vegetarianism, it is more likely that he was influenced by Hindu tradition. Even his mathematics has an Indian flavour.

About 525 BC, Pythagoras emigrated to Croton (modern Crotone) in southern Italy, where he founded a society, half-way between a political party and a religious cult, which came to be known as the ‘Pythagorean Brotherhood.’ Some members of this society were admitted to an inner circle consisting of the so-called ‘mathematicians’.

The word ‘mathematics’ was in fact introduced by Pythagoras. The first part of this word is an old Indo-European root, related to the English word ‘mind’. The modern meaning of ‘mathematics’ is due to Aristotle.

Whereas Thales had claimed that ‘all is water’, Pythagoras taught that ‘all is number’. For Pythagoras this implied that everything could be understood in terms of whole numbers and their ratios. In particular, he implicitly expected that every line segment was a whole number or a ratio of whole numbers (in terms of a given unit length). It seems that the dis-

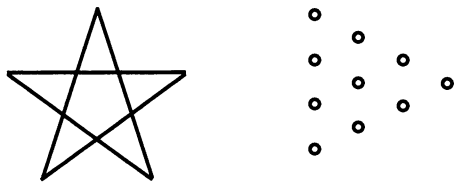


FIGURE 7.1. Pythagorean star and the fourth triangular number

covery of the irrationality of the diagonal of the square of side 1 was made by his followers and that Pythagoras himself was not aware of this.

In his philosophy, Pythagoras reserved a special place for the number 10. He called it the ‘divine number’, noting that 10 is a triangular number and realizing that the five-pointed ‘Pythagorean star’ (Figure 7.1) has 10 vertices.

The Pythagoreans ascribed all their mathematical discoveries to Pythagoras, but there is not, in fact, a single theorem which we can safely credit to him. For example, in his preface to the *Introductio Arithmetica*, written by a Pythagorean, Nichomachus of Gerasa (100 AD), Iamblichus (300 AD) credits Pythagoras with a knowledge of the amicable pair 220 and 284. (Two natural numbers are *amicable* if each is the sum of the proper divisors of the other.) However, we have no way of knowing for certain whether amicable numbers had been recognized as early as 500 BC. Yet, according to a famous anecdote, when someone challenged his slogan ‘all is number’ by asking ‘then what is friendship?’, Pythagoras replied that friendship is as 220 is to 284.

Leaving behind the shadowy figure of the Master, let us review the accomplishments of his followers. Although they were primarily a religious and political group, they did a fair amount of work in arithmetic, geometry, astronomy and music – the four subjects later forming the medieval *quadrivium*. (In the university curriculum of the Middle Ages, these subjects were meant to follow the ‘trivial’ subjects: grammar, rhetoric and logic.)

Theorem of Pythagoras

The Pythagoreans are probably responsible for the proof found in Euclid’s *Elements*, Book I, Proposition 47. They also found a proof of the converse of the theorem of Pythagoras.

Means

They examined the relationships between the following means:

arithmetic ($\frac{1}{2}(a + b)$), geometric (\sqrt{ab}) and harmonic ($2ab/(a + b)$).

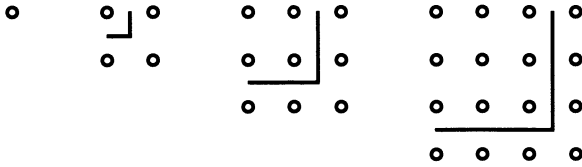


FIGURE 7.2. The sequence of squares

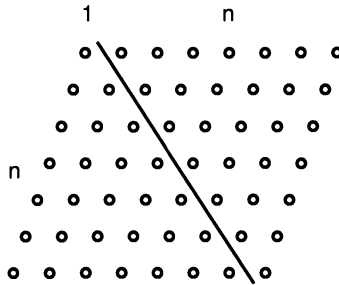


FIGURE 7.3. Triangular numbers

Perfect Numbers

They found a formula giving perfect numbers. See Chapter 8.

Regular Solids

They discovered the dodecahedron. See Chapter 9.

Irrationality of $\sqrt{2}$

They discovered that the square root of 2 is not rational. They used the integer solutions of $x^2 - 2y^2 = \pm 1$ to find approximations to it. See Chapter 10.

Figurative Numbers

They found proofs for several algebraic relations by means of studying figurative numbers. For example, looking at the sequence of squares, expressed in terms of ‘arrays of pebbles’ (Figure 7.2), they noticed that $n^2 + (2n + 1) = (n + 1)^2$ and hence $1 + 3 + 5 + \cdots + (2n - 1) = n^2$:

Fitting two ‘triangular numbers’ into a parallelogram, they noticed that the n th triangular number is $\frac{1}{2}n(n + 1)$.

Looking at the sequence of triangular numbers, expressed in terms of pebble arrays, they realized that the difference between the $(n + 1)$ th and n th triangular number is just $n + 1$. From this they concluded that $1 + 2 + 3 + \cdots + n =$ the n th triangular number $= \frac{1}{2}n(n + 1)$. See Figure 7.3.

The study of figurative numbers is alive and well today. For example, recently some very advanced mathematics was used to show, for the first

time, that there are exactly six triangular numbers that are products of three consecutive integers (See Tzanakis and de Weger [1989].)

Exercises

1. How might a Pythagorean have derived the fact that the angle at the tip of his star is 36° ?
2. Check that 220 and 284 are amicable. If 12,285 is one member of an amicable pair, find the other.
3. Prove that the sum of the first n cubes is the square of the sum of the first n numbers. (Use mathematical induction or construct a square figure with sides of length $1 + 2 + \cdots + n$ and divide it into figures whose areas are the first n perfect cubes.)

8

Perfect Numbers

The Pythagoreans were interested in perfect numbers, that is, numbers such as 6 and 28, which are equal to the sum of their proper divisors. They may also be described as numbers which are amicable with themselves. Nowadays we usually speak about the sum of all the divisors of a positive integer n , including n itself. If $\sigma(n)$ denotes this sum, then n is *perfect* if and only if $\sigma(n) = 2n$. As the culmination of Book IX of the *Elements* (300 BC), Euclid proved that any positive integer of the form

$$n = 2^{m-1}(2^m - 1)$$

is perfect, whenever $2^m - 1$ is prime. This fact had probably been discovered by the Pythagoreans.

Proof of Proposition IX 36 (Perfect Number Theorem):

If $p = 2^m - 1$ is prime, then the divisors of $n = 2^{m-1}p$ are

$$1, 2, 2^2, \dots, 2^{m-1}, p, 2p, 2^2p, \dots, 2^{m-1}p$$

(thanks to unique factorization). The sum of these divisors is thus

$$\begin{aligned}\sigma(n) &= (1 + 2 + 2^2 + \dots + 2^{m-1})(1 + p) \\ &= (2^m - 1)(1 + p) \\ &= 2(2^{m-1}(2^m - 1)) = 2n.\end{aligned}$$

Even though $2^m p$ is not square-free, Euclid did have a rigorous proof of the special case of the unique factorization theorem which is used in the

2	107	9689
3	127	9941
5	521	11213
7	607	19937
13	1279	21701
17	2203	23209
19	2281	44497
31	3217	86243
61	4253	132049
89	4423	216091

TABLE 8.1. Values of m making $2^m - 1$ prime

above proof, and he also had a rigorous proof for the formula for the sum of a geometric progression (IX 35).

An integer of the form $2^m - 1$ can only be prime if m is prime. For if $m = ab$ with $a, b > 1$, we have the factorization

$$2^{ab} - 1 = (2^a - 1)((2^a)^{b-1} + (2^a)^{b-2} + \cdots + 2^a + 1)$$

into two factors greater than 1. The converse is not true. Although 11 is prime, $2^{11} - 1$ is not; for $2^{11} - 1 = 2047 = 23 \times 89$.

Primes of the form $2^m - 1$ are called *Mersenne* primes, after Father Marin Mersenne (1588–1648). In the preface of his *Cogitata Physico-Mathematica* (1644), Mersenne correctly stated that the first 8 perfect numbers are given by the values $m = 2, 3, 5, 7, 13, 17, 19$ and 31. He also claimed that $2^{67} - 1$ is prime, and hence $2^{66}(2^{67} - 1)$ is perfect. Here he was wrong. In 1903, Frank Nelson Cole gave a lecture which consisted of two calculations. First Cole calculated $2^{67} - 1$. Then he worked out the product

$$193, 707, 721 \times 761, 838, 257, 287.$$

He did not say a word as he did this. The two calculations agreed, and Cole received a standing ovation. He had factored $2^{67} - 1$ and proved Mersenne wrong.

Edouard Lucas (1842–1891) found a very efficient way of testing whether $2^m - 1$ is prime. Let $u_1 = 4$ and $u_{n+1} = u_n^2 - 2$. Thus $u_2 = 14, u_3 = 194$ and $u_4 = 37,634$. If $m > 2$ then $2^m - 1$ is prime just in case $2^m - 1$ is a factor of u_{m-1} . For example, since $2^5 - 1$ is a factor of 37,634 it follows that $2^5 - 1$ is prime (and hence $2^4(2^5 - 1) = 496$ is perfect). (For an elementary proof of Lucas's Theorem, see Sierpinski [1964].)

Thanks to Lucas's test — and the modern computer — we know, since about 1985, that $2^m - 1$ is prime when m has the 30 values given in Table

8.1. The Greeks knew just the first four Mersenne primes and Mersenne discovered eight more. Before 1950, only the first 12 Mersenne primes were known. Then, with the help of ever more powerful computers, 18 more came to light. (Even after writing this, we learned of three more, corresponding to $m = 110,503$, $m = 756,839$, and $m = 858,433$. The last of these Mersenne primes has 258,716 digits.) We still do not know whether there are infinitely many Mersenne primes. Nor do we know if there are any odd perfect numbers. What we do know is that every even perfect number has the form given by Euclid. This was first proved by Leonhard Euler (1707–1783). His proof was as follows.

Proof that every even perfect number has Euclid's form:

Suppose n is an even perfect number. Then we can write it in the form $2^{m-1}q$ with q odd and $m, q > 1$. Each divisor of n has the form $2^r d$ where $0 \leq r \leq m-1$ and d is a divisor of q . Therefore

$$\sigma(n) = (1 + 2 + \dots + 2^{m-1})\sigma(q) = (2^m - 1)\sigma(q).$$

Since n is perfect, $2^m q = \sigma(n) = (2^m - 1)\sigma(q)$.

Since $2^m - 1$ is odd, 2^m must divide $\sigma(q)$, say $\sigma(q) = 2^m k$, hence $q = (2^m - 1)k$. Among the divisors of q are q itself and k . These are different, since $m > 1$, and their sum is $2^m k$, which is the sum of all the divisors of q . Therefore, q has exactly two divisors and so is prime, hence $k = 1$ and $q = 2^m - 1$.

Perfect numbers are of interest not only as a challenge to computer programmers, they also play role in religious mysticism. For example, following Philo of Alexandria, Augustine writes in the *City of God*:

Six is a number perfect in itself, and not because God created all things in six days; rather, the converse is true. God created all things in six days because this number is perfect, and it would have been perfect even if the work of six days did not exist.

In a recent book on Sufi mysticism, it is stated that 6 is the first 'complete' number and 28 is the second. Evidently, 'complete' here means 'perfect'.

Apparently, the Pythagoreans knew only one amicable pair of numbers. Although Euler found 60 such pairs, the second smallest pair (1184, 1210) was only discovered in 1866 by Nicolo Paganini.

The Arabic mathematician Thabit ibn Qurra (826–901) gave a general procedure for discovering many amicable pairs, analogous to Euclid's procedure for discovering perfect numbers. See Exercise 5.

Exercises

1. Prove that every even perfect number, except 6, is the sum of the first 2^k odd cubes, for some k .
2. Show that, if m and n are relatively prime positive integers, then $\sigma(mn) = \sigma(m)\sigma(n)$.
3. Show that, if p is prime, $\sigma(p^k) = (p^{k+1} - 1)/(p - 1)$.
4. Obtain a formula for $\sigma(n)$ in terms of the prime factorizations of n .
5. Prove the result of Thabit ibn Qurra: if $p = 3 \times 2^{t+1} - 1$, $q = 3 \times 2^t - 1$ and $r = 9 \times 2^{2t+1} - 1$ are odd prime numbers, then $m = 2^{t+1}pq$ and $n = 2^{t+1}r$ are amicable.
6. Find two amicable pairs with the help of the above procedure.

9

Regular Polyhedra

The Pythagoreans knew that there are three ways to tile a plane (e.g., a bathroom floor) using congruent regular polygons. Indeed, since once can dissect a polygon with p sides into $p - 2$ triangles, the sum of the angles of such a polygon is $(p - 2)180^\circ$. Thus each angle of a regular ' p -gon' is $(p - 2)180^\circ/p$. If q such angles meet at a point, then

$$q(p - 2)180^\circ/p = 360^\circ ,$$

which may be simplified to yield the Egyptian problem: $1/2 = 1/p + 1/q$. Since p and q are integers greater than 2, we must have one of the following three possibilities:

p	q
3	6
4	4
6	3

The first of these gives the tiling with equilateral triangles, the second the tiling with squares, and the last the tiling with regular hexagons. No other regular polygon can be used to tile the plane.

A polyhedron is *regular* if its faces are congruent regular polygons and if the same number of faces meet at each vertex. Five regular polyhedra are the following:

- the cube, bounded by 6 squares, with 3 edges at each vertex,
- the tetrahedron, bounded by 4 equilateral triangles, with 3 edges meeting at each vertex,

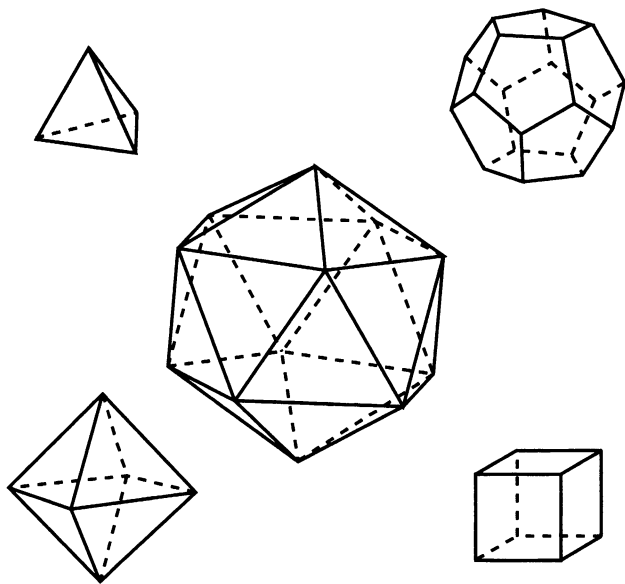


FIGURE 9.1. The Pythagorean solids

- the octahedron, bounded by 8 equilateral triangles, with 4 edges at each vertex,
- the icosahedron, bounded by 20 equilateral triangles, with 5 edges at each vertex,
- the dodecahedron, bounded by 12 regular pentagons, with 3 edges at each vertex (see Figure 9.1).

In the *Timaeus* (53-58), Plato explains the composition of the physical universe in terms of these five regular polyhedra. The cube is associated with earth, the tetrahedron with fire, the octahedron with air, the icosahedron with water and the dodecahedron with the whole cosmos. Plato explains the boiling of water by means of a ‘chemical equation’, which we might write as

$$F_4 + W_{20} \rightarrow 2A_8 + 2F_4.$$

That is, fire, with 4 faces, combines with water (20 faces) to produce 2 air atoms (each with 8 faces) and 2 fire atoms (each with 4 faces). Note that the numbers balance.

Such theorizing is very much in the spirit of the Pythagorean teaching that ‘all is number’. Indeed, historians attribute the theory of the *Timaeus* to Pythagoras himself (Guthrie [1987]). It seems, however, that Pythagoras himself may not have known about the dodecahedron. According to one account, Hippasus (470 BC) was expelled from the Pythagorean or-

p	q	$1/p + 1/q$	E	F	V	Polyhedron
3	3	$2/3$	6	4	4	tetrahedron
3	4	$7/12$	12	8	6	octahedron
4	3	$7/12$	12	6	8	cube
3	5	$8/15$	30	20	12	icosahedron
5	3	$8/15$	30	12	20	dodecahedron

TABLE 9.1. The regular polyhedra

der because, having discovered the dodecahedron, he failed to ascribe his discovery to Pythagoras.

When Plato proposed that the creator, whom he called the ‘demiurge’, used the regular polyhedra when forming the universe, he may not have been that far off. The tetrahedron, cube and octahedron can be found in nature as crystals. The octahedron, icosahedron and dodecahedron occur as the skeletons of certain radiolarians (a type of microscopic sea animal).

Are there only five regular polyhedra? Yes. In fact, a proof of this is found at the end of Euclid’s *Elements* (300 BC). This proof is based on the fact that if q regular p -gons meet at a vertex, then the sum of the q angles in the q faces is less than 360° . This is proved rigorously in Proposition 21 of Book XI of the *Elements*, but it can be seen intuitively by imagining someone cutting the q edges and flattening the angle. For example, the three angles at the vertex of a cube clearly add up to 270° , which is less than 360° .

In general, for a regular polyhedron whose faces are regular p -gons, with q faces meeting at each vertex,

$$q(p-2)180^\circ/p < 360^\circ,$$

which may be simplified to yield $1/2 < 1/p + 1/q$. We can easily see that there are just five possibilities for p and q , as in Table 9.1. (In the table, E is the number of edges the polyhedron has, F the number of faces, and V the number of vertices.)

We have not as yet explained how the numbers E , F and V in our table are calculated. It so happens that these numbers are related by a simple formula, even when we consider an arbitrary polyhedron, regular or not. Here are some examples:

	F	V	E
cube	6	8	12
tetrahedron	4	4	6
pyramid	5	5	8
prism	5	6	9

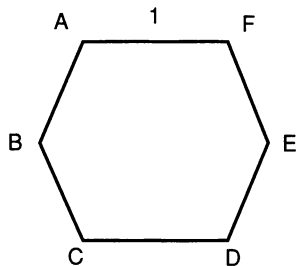


FIGURE 9.2. Cross-section of an icosahedron

We note that in each example

$$F + V - E = 2.$$

This is in fact a general rule, valid for all polyhedra. While it may first have been noted by Descartes, it was only proved by Euler and is known as *Euler's formula*. We shall give a proof in Chapter 30 of Part I.

Now suppose we are looking at a regular solid in which each face has p edges and in which q edges meet at each vertex. It follows immediately that

$$pF = 2E, \quad qV = 2E.$$

Substituting $F = 2E/p$ and $V = 2E/q$ into Euler's formula, we obtain

$$2E/p + 2E/q - E = 2,$$

and, after dividing by $2E$,

$$1/p + 1/q - 1/2 = 1/E.$$

This allows us to calculate E from p and q , and then F and V .

The ancient Greeks were fascinated by the five regular solids. Without the help of trigonometry or calculus, they managed to prove all their basic properties. Book XIII of the *Elements* (300 BC) is devoted to showing that, for each of these five solids, there is a sphere passing through all its vertices. In each of the five cases, Book XIII calculates the ratio of the side of the regular polyhedron to the radius of this 'circumscribing' sphere.

For example, if one cuts an icosahedron in half, cutting along a side, the resulting cross-section is as in Figure 9.2. AF and CD are edges; AC and DF are diagonals in the regular pentagons formed by the sides of the icosahedron. (You may have to construct an icosahedron to see this.) Thus, if AF and CD each have unit length, AC and DF each have length $\frac{1}{2}(1 + \sqrt{5})$. The diameter of the circumscribing sphere is CF , which is the hypotenuse of the right triangle with sides CD and DF . Thus $CF^2 = 1^2 + (\frac{1}{2}(1 + \sqrt{5}))^2$, and hence the radius of the sphere is $\frac{1}{2}\sqrt{\frac{1}{2}(5 + \sqrt{5})}$.

Exercises

1. Construct a dodecahedron, e.g., by taping together 12 identical regular pentagons cut out of cardboard.
2. Show that the radius of a sphere passing through the vertices of a dodecahedron with side 1 is $(\sqrt{3} + \sqrt{15})/4$.
3. Show that the volume of a dodecahedron of side 1 is $(15 + 7\sqrt{5})/4$.
4. Given a polyhedron, not necessarily regular, in which exactly 3 edges meet at each vertex. Show that $V = 2K$, $E = 3K$ and $F = K + 2$ for some positive integer K .
5. Under the condition of the previous exercise, if F_p is the number of faces with p sides, show that

$$\sum_p (6 - p)F_p = 12.$$

6. If all faces of a polyhedron are hexagons or pentagons and if three edges meet at each vertex, prove that the number of pentagons is twelve. (There are molecules of such a shape with twenty hexagons, called 'buckyballs', a form of carbon called 'buckminster fullerene'. See Chung and Sternberg [1993].)

The Crisis of Incommensurables

Two lengths a and b are said to be *commensurable* if there exist positive integers p and q such that $a/b = p/q$. When the Pythagoreans claimed that all things are numbers, they probably meant to imply that all pairs of lengths are commensurable. They were aware of the fact that, if a vibrating string is divided into two parts, of lengths a and b , so that a melodious tone is produced, then a and b are commensurable.

Unfortunately for the Pythagoreans, they soon discovered that the diagonal of a square is not commensurable with its side. A simple proof of this is found in Aristotle's *Prior Analytics* 41a23-30. Let $ABCD$ be a square, say of side $AB = 1$. By the Theorem of Pythagoras, the diagonal AC measures $\sqrt{2}$. Suppose $\sqrt{2} = AC/AB = p/q$, where p and q are positive integers. We may assume, without loss of generality, that p and q are relatively prime. In particular, they are not both even. Now $p^2 = 2q^2$, so that p^2 is even. As the Pythagoreans well knew, the square of an odd number is odd and the square of an even number is even. Thus, from the fact that p^2 is even, it follows that p is even. Putting $p = 2r$, we have $2q^2 = (2r)^2$, hence $q = 2r^2$. But this means that q is even as well, contradicting the fact that p and q are relatively prime. Thus, the assumption that AC and AB are commensurable must be false. Today we would express this result by saying that $\sqrt{2}$ is irrational.

The Pythagoreans tried to keep this discovery a secret, as it seemed to undermine their whole philosophy. Some say that it was Hippasus, whom we met before, who leaked the secret, and that he drowned in a shipwreck as a punishment for having done so. It seems that Hippasus was the Trotsky of the Pythagorean society.

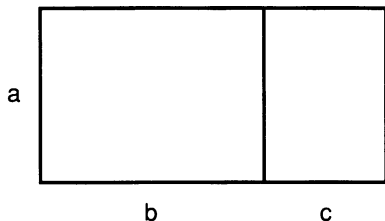
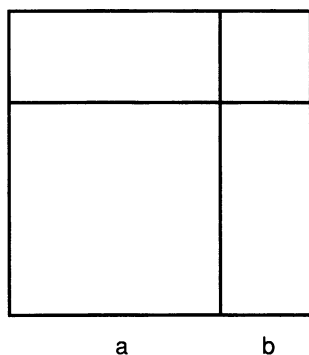


FIGURE 10.1. The distributive law

FIGURE 10.2. Binomial expansion $(a + b)^2 = a^2 + 2ab + b^2$

The Greeks did not have infinite decimals. They did not know how to handle a number like $\sqrt{2}$ in an arithmetical or algebraic fashion the way we do now, although it has recently been claimed that they could represent real numbers by continued fractions. They did, however, know that $\sqrt{2}$ was a length, and they turned to geometry for an understanding of it. The problem of incommensurables was one of the reasons that they preferred to do what we would call algebra in a geometric manner.

For example, the distributive law $a(b + c) = ab + ac$ was thought of as an addition rule for areas of rectangles, as in Figure 10.1.

Euclid put it thus:

If there are two straight lines, and one of them be cut into any number of segments whatever, the rectangle contained by the two straight lines is equal to the rectangles contained by the uncut straight line and each of the segments (*Elements* II 1).

The law $(a + b)^2 = a^2 + 2ab + b^2$ is illustrated in Figure 10.2. We shall refrain from putting this law into words.

However, as a third example, we again quote Euclid:

If a straight line be cut into equal and unequal segments, the rectangle contained by the unequal segments of the whole to-

gether with the square on the straight line between the points of section is equal to the square on the half (*Elements* II 5).

This is equivalent to the identity $(a + b)(a - b) = a^2 - b^2$.

The Pythagoreans found ways of approximating $\sqrt{2}$ as closely as could be desired by rational numbers. Using our modern algebraic notation, we can express their method as follows.

If $x^2 - 2y^2 = \pm 1$, with x and y positive integers, then x^2 is approximately equal to $2y^2$, so that x/y is approximately equal to $\sqrt{2}$. More precisely,

$$(x/y - \sqrt{2})(x/y + \sqrt{2}) = x^2/y^2 - 2 = \pm 1/y^2$$

so that

$$x/y - \sqrt{2} = \pm 1/(y^2(x/y + \sqrt{2})).$$

Since $x/y + \sqrt{2} > 1$, it follows that

$$|x/y - \sqrt{2}| < 1/y^2.$$

Thus, if we can find positive integer solutions of $x^2 - 2y^2 = \pm 1$ with y sufficiently large, then we can find rational approximations to $\sqrt{2}$ as close as we please.

To find positive integers x and y such that $x^2 - 2y^2 = \pm 1$, the Pythagoreans proceeded as follows. Putting

$$a_1 = 1, \quad b_1 = 1$$

and defining inductively

$$a_{n+1} = a_n + 2b_n, \quad b_{n+1} = a_n + b_n,$$

they obtained the following table:

n	a_n	b_n	a_n/b_n
1	1	1	1
2	3	2	3/2
3	7	5	7/5
4	17	12	17/12

etc., in which the last column contains successive approximations to $\sqrt{2}$.

Indeed, it is not difficult to prove by mathematical induction that

$$a_n^2 - 2b_n^2 = (-1)^n.$$

This is surely true when $n = 1$, so suppose it holds for n . Then

$$\begin{aligned}
 a_{n+1}^2 - 2b_{n+1}^2 &= (a_n + 2b_n)^2 - 2(a_n + b_n)^2 \\
 &= a_n^2 + 4a_nb_n + 4b_n^2 - 2a_n^2 - 4a_nb_n - 2b_n^2 \\
 &= -a_n^2 + 2b_n^2 \\
 &= -(-1)^n = (-1)^{n+1}.
 \end{aligned}$$

Thus our proof by mathematical induction is complete.

This method, but not the above proof, is explained verbally by Proclus in commenting on a passage in Plato's *Republic*. Today, it is easy to obtain explicit formulas for the numbers a_n and b_n . First, one proves by mathematical induction that

$$a_n + b_n\sqrt{2} = (1 + \sqrt{2})^n .$$

Replacing the square root by its negative, one obtains

$$a_n - b_n\sqrt{2} = (1 - \sqrt{2})^n .$$

Therefore,

$$\begin{aligned} a_n &= \frac{1}{2}((1 + \sqrt{2})^n + (1 - \sqrt{2})^n) , \\ b_n &= \frac{1}{2\sqrt{2}}((1 + \sqrt{2})^n - (1 - \sqrt{2})^n) . \end{aligned}$$

Although the Pythagoreans did not know it, they had actually found all solutions of the equations $x^2 - 2y^2 = \pm 1$ in positive integers. Suppose, for example, $x^2 - 2y^2 = 1$. Let n be the largest natural number such that $(1 + \sqrt{2})^n \leq x + y\sqrt{2}$, then

$$(1 + \sqrt{2})^n \leq x + y\sqrt{2} < (1 + \sqrt{2})^{n+1} .$$

Multiplying this by $(1 - \sqrt{2})^n = a_n - b_n\sqrt{2}$ and assuming that n is even, we obtain

$$(1) \quad 1 \leq (x + y\sqrt{2})(a_n - b_n\sqrt{2}) < 1 + \sqrt{2} .$$

Taking reciprocals of this, we get

$$(2) \quad -1 \leq (-x + y\sqrt{2})(a_n + b_n\sqrt{2}) < 1 - \sqrt{2} .$$

Adding (1) and (2) and dividing by $2\sqrt{2}$, we obtain

$$0 \leq ya_n - xb_n < 1/\sqrt{2} .$$

Since $ya_n - xb_n$ is a whole number, it must be 0, hence $ya_n = xb_n$. Now we know that x and y are relatively prime, and so are a_n and b_n . It easily follows that $x = a_n$ and $y = b_n$, where n is even.

If n is odd or if $x^2 - 2y^2 = -1$, we proceed similarly.

Exercises

1. Prove that the decimal expression of $\sqrt{2}$ is not ultimately periodic.
2. Prove that the following numbers are not rational: $\sqrt{3}$, $\sqrt[3]{2}$ and $\log_{10} 2$.
3. If a, b, c and d are integers and $a + b\sqrt{2} = c + d\sqrt{2}$, show that $a = c$ and $b = d$.
4. Solve the following equations for positive integers:

$$x^2 - 4y^2 = 1, \quad x^2 - 3y^2 = 1.$$

From Heraclitus to Democritus

Heraclitus of Ephesus (in western Turkey) flourished about 500 BC, Parmenides of Elea (in southern Italy) about 480 BC, Zeno of Elea about 460 BC, Empedocles in Sicily about 440 BC, Democritus of Abdera (in north-eastern Greece) about 420 BC.

In the *Metaphysics* (986b4-8), Aristotle tells us that the Pythagoreans had a list of opposites: one, many; finite, infinite; male, female; etc. It was perhaps this list which led Heraclitus to his view that everything that happens is the result of a struggle between opposites. He proclaimed that all change is the result of strife.

Heraclitus believed that everything is in flux. It was he who asserted that one cannot step into the same river twice. Not surprisingly, he thought the fundamental substance was fire, and declared that all matter can be transformed into fire (and vice versa), just as all goods can be exchanged for gold. Did he anticipate the modern discovery that mass can be transformed into energy?

Heraclitus has had a great deal of influence on the twentieth century, largely through the nineteenth century Prussian philosopher Hegel. Influenced by Heraclitus, Hegel taught that the universe is a sort of debating society in which 'thesis' and 'antithesis' are forever struggling to produce a 'synthesis'. Marx adopted this philosophy, giving it a materialistic slant, and the views of Heraclitus ended up forming part of the official doctrine of Marxist governments, now much in decline.

Heraclitus has had less influence on logic. On one occasion he expressed his doctrine of continual change by saying that the river we step into both is and is not the same. Yet, in most logical systems, any statement of the

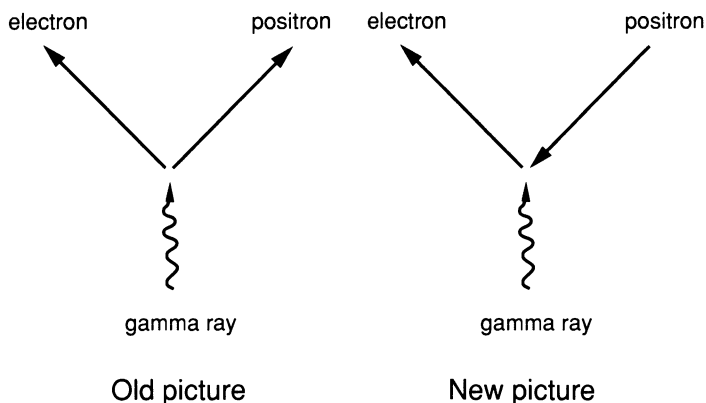


FIGURE 11.1. Positron as an electron travelling backwards in time

form ' p and not p ' is regarded to be false. Hegelians sometimes adopt a similar mode of speech, claiming that ' a is not always equal to a '. Needless to say, this doctrine has not been applied to mathematics.

Yet Marxist philosophers try to understand not only history, but also mathematics in terms of a dialectic process. According to Lenin, subtraction is the antithesis of addition, yielding arithmetic as a synthesis, and integration is the antithesis of differentiation, the synthesis being calculus. Quite recently, the American mathematician Lawvere has suggested that a foundation of mathematics be built on a dialectic process in which the striving opposites are so-called 'adjoint functors', but this concept is too technical to be explained here.

Parmenides took the view opposite to that of Heraclitus, proclaiming that nothing changes, that change is an illusion: from the point of view of the 'goddess', the past and the future are all there at the same time. This is a bit like the view of the modern physicist and his four-dimensional space-time, in which the ever-changing events are replaced by unchanging world-lines.

Richard Feynman has recently shown that this way of viewing the universe allows one to give a more elegant and instructive explanation of certain fundamental processes. For example, to explain how an electron and a positron annihilate each other, giving rise to a γ ray, we may take it that the positron is an electron traveling backwards in time, having been deflected with a γ ray splitting off. Simultaneous pair creation is explained similarly. One may even speculate that there is only one electron in the universe. See Figure 11.1.

Zeno was a disciple of Parmenides. He produced four arguments attempting to prove that motion is impossible, his so-called 'paradoxes'. What he

really showed was that if you do not allow infinite processes, what we now call 'limits' in mathematics, then you cannot use mathematics to analyze motion. These arguments, found in Aristotle's *Physics* 239b5-240a18 and 233a21-31, are the following.

1. A point moving from 0 to 1 on the number line first covers a distance of $1/2$, then a distance of $1/4$, then a distance of $1/8$, and so on. After n steps, it has covered a total distance of $1/2 + 1/4 + \dots + 1/2^n = 1 - 1/2^n$. From the fact that there is no n such that $1 - 1/2^n = 1$, Zeno concluded that the point will never reach 1. In other words, motion from 0 to 1 is impossible.

Today we get around Zeno's difficulty with the help of the notion of 'limit'. But this is a fairly sophisticated concept; according to nineteenth century mathematicians, the meaning of ' $\lim_{n \rightarrow \infty} f(x) = a$ ' is as follows: for every real $\epsilon > 0$, there exists a natural number k such that, for all $n > k$, $|f(n) - a| < \epsilon$. While the ancients may have had an intuitive notion of what is meant by a limit, the rigorous definition was surely beyond them.

2. Achilles and the tortoise are engaged in a race along a measured line. Achilles starts at 0, but the tortoise is given a head start, at 1, since Achilles runs twice as fast as the tortoise. Thus, when Achilles arrives at the point 1, the tortoise is at $1 + 1/2$; when Achilles arrives at the point $1 + 1/2$, the tortoise is at $1 + 1/2 + 1/4$, and so on. In general when Achilles gets to $2 - 1/2^{n-1}$, the tortoise is at $2 - 1/2^n$, just a little bit ahead. From this Zeno concludes that Achilles can never catch up with the tortoise. If it looks like Achilles catches up with the tortoise, that only means that motion is an illusion.

The modern solution to this paradox is the same as the solution to paradox (1): $\lim_{n \rightarrow \infty} (2 - 1/2^{n-1}) = \lim_{n \rightarrow \infty} (2 - 1/2^n)$.

3. Since, at any instant, a flying arrow is in exactly one place, Zeno infers that, at that instant, it is motionless. (Indeed, if you took a high-speed photograph of the arrow, it would look as if it was perfectly still.) One is tempted to conclude that the speed of the arrow is 0, that it is not really moving.

The argument seems to be that the speed $dx/dt = 0$ because $dx = 0$. But this would only follow if $dt \neq 0$. If we assume that an interval of time is an actually infinite set of (equal) instants, then each instant has zero duration, and $dt = 0$. But then we could equally well infer that $dx/dt = 17$, since $0 = 17 \times 0$.

Zeno's argument works only if you assume that time is basically discrete, there being some smallest, finite 'quantum' of time (e.g., $1/2^{100}$ seconds). During each quantum of time, the arrow would be motionless, for if it moved, say, from 0 to 1, there would be a time before

it got to $1/2$, and a time after, and the quantum of time would be divisible into two parts. Hence it really would not be moving. If there were some n such that every interval of time exceeds $1/2^n$ seconds, then Zeno would be right: motion is an illusion.

In writing the speed as dx/dt , we used the notation of the 17th century philosopher and mathematician Leibniz. He conceived of dx and dt as *infinitesimals* and thought of dx/dt as their actual ratio. The great Newton essentially shared this view, even though his notation was different. Infinitesimals were believed to be quantities which are infinitely small, yet unequal to zero. This idea was attacked by the 18th century philosopher Berkeley as being absurd, and 19th century mathematicians agreed with him. They redefined the ratio dx/dt as

$$\frac{dx}{dt} = \lim_{\delta t \rightarrow 0} \frac{\delta x}{\delta t},$$

where δx and δt are not necessarily small.

It was only in the middle of the 20th century that Abraham Robinson pointed out that infinitesimals may be introduced by fiat, just like the square root of -1 , as *their existence does not lead to a contradiction*. Indeed, consider the following infinite collection of inequalities:

$$0 < dx, \quad dx < 1, \quad dx < 1/2, \quad dx < 1/3, \quad \dots \quad (*)$$

Suppose we can derive a contradiction from this infinite collection. Now, it is generally agreed that a mathematical proof can have only a finite number of steps, this being part of the very definition of *proof*. Therefore, the proof of a contradiction from the assumptions $(*)$ can only mention a finite number of them, say the last being $dx < 1/n$. But this finite collection of assumptions does not lead to a contradiction, as $dx = 1/(n+1)$ satisfies all of them.

4. In Zeno's fourth argument there are three rows of people:

$$\begin{array}{ccccccc} & A & A & A & A & & \\ & B & B & B & B & \rightarrow & \\ \leftarrow & C & C & C & C & & \end{array}$$

The A 's are stationary, the B 's are moving to the right and the C 's are moving to the left, at the same speed. Now suppose it takes a B one instant of time to pass an A , then it will take him half an instant of time to pass a C . It follows, that there is no such thing as an *instant*, if by that we mean an indivisible 'quantum of time': time is infinitely divisible, it is a substance.

In summary, we may not agree with Zeno that there is no motion, but we must credit him with probing into the very foundation of mathematics and physics.

Empedocles is important in experimental science and cosmology. He demonstrated that air is a substance by pushing an inverted tumbler into a tub of water; the water did not rush in to fill the apparent vacuum. He recognized not only the four traditional substances, earth, water, air and fire, but postulated two other: love and strife. He believed that, as long as love prevails, the cosmos contracts, but when strife prevails, it expands. Thus he seems to have anticipated modern astronomers in realizing that our universe is in an expanding phase. As his last experiment, he leaped into a volcano to demonstrate his immortality.

Democritus finally did away with substances and replaced them with atoms. These were assumed to be physically but not geometrically, indivisible. They were indestructible and constantly moving, the space between them being empty. The number of atoms was infinite, but they differed in shape and size. In the 19th century, people thought they had identified the atoms of Democritus; but then, in the early 20th century, Rutherford showed that the so-called atoms were divisible after all. It might have been wiser to reserve the word 'atom' for our electrons and quarks, if these indeed turn out to be the ultimate constituents of matter.

Democritus was a determinist, he believed that everything that happens must happen. He wrote books on geometry, which have been lost. He is said to have emphasized the importance of proofs and to have discovered (or rediscovered) how to calculate the volume of a pyramid or cone by taking one third of the base area times the height. It is curious that Marx wrote his doctoral dissertation on Democritus and not on Heraclitus.

12

Mathematics in Athens

After an alliance of Greek cities defeated the Persians in 490 BC, Athens became, for over a hundred years, a great center of civilization. There are things about it we do not like: Athens exacted tribute from its allies, and the leisurely life of its leading citizens was based upon slavery. Nonetheless, one can safely say that the degree of civilization achieved by the Athenians around 400 BC has rarely, if ever, been surpassed in the history of the world. Because we must confine our attention to mathematics, we shall touch on only one of the many areas in which cultural development took place.

One of the first mathematicians who worked in Athens was the sophist Hippias (420 BC), from Elis on the west coast of Greece. In a dialogue sometimes ascribed to Plato, we hear Socrates (469–399 BC) teasing Hippias about his mathematics:

Socrates: And tell me, Hippias, are you not a skilful calculator and arithmetician?

Hippias: Yes, Socrates, assuredly I am.

Socrates: And if someone were to ask you what is the sum of 3 multiplied by 700, you would tell him the true answer in a moment, if you pleased?

Hippias: Certainly I should.

Socrates: Is not that because you are the wisest and ablest of men in these matters?

Hippias: Yes.

Hippias discovered a curve called the *quadratrix*, which can be used for trisecting an arbitrary angle, and also for constructing a square equal in area to a given circle. He described the quadratrix as follows: imagine that side AD of the unit square $ABCD$ moves down at a rate of 1 unit per second towards the side BC (on the ‘bottom’ of the square). Imagine that side AB rotates about B at a rate of $1/4$ revolution per second towards BC , so that, after 1 second, both AD and AB coincide with BC . At any time t ($0 \leq t \leq 1$), the two moving sides meet at a point P . The set or *locus* of these points P is the quadratrix.

In terms of our modern analytic geometry and trigonometry, we would put it this way: the point P has coordinates

$$((1-t)/\tan(90^\circ(1-t)), 1-t)$$

so the equation of the quadratrix is $y = x \tan(90^\circ y)$, with $0 \leq y$.

To divide an angle of, say, 60° into 3 equal parts, it is enough to place it in standard position, with its vertex on the origin, and one arm along the positive x axis. If the other arm meets the quadratrix at (a, b) , we find the point P where $y = b/3$ meets the quadratrix. The angle between the line joining P to the origin and the positive x axis is 20° .

Furthermore, as y tends to 0, $x = y/\tan(90^\circ y)$ tends to $2/\pi$, and so the quadratrix can be used to ‘square the circle’.

This highly ingenious method was criticized — by Plato, it seems — on the grounds that it is more elegant to use only straight lines and circles in the solution of mathematical problems. One ought to carry out geometric constructions using only a ruler (for drawing straight lines) and a compass (for drawing circles); using a quadratrix was considered to be cheating.

However, as Pierre Wantzel (1814–1848) was the first to prove, it is not possible to trisect an arbitrary angle using only straight lines and circles. One has to use some other tool — such as the quadratrix. Hippias was right and Plato was wrong; although by insisting on ruler and compass constructions he raised an interesting and challenging problem, which we shall discuss in Chapter 14.

In order to understand the Greek contribution to the beginnings of analysis, it is important to know how they attacked, and finally solved, the problem of the area of the circle. Antiphon the sophist (425 BC) was one of the early Athenian mathematicians who worked on this problem. He suggested that the area of the circle be calculated in terms of the regular polygons inscribed in it. (The *regular m-gon* is the polygon with m equal sides and angles.)

Using the assumption that the area of the union of pairwise disjoint sets equals the sum of their areas, it is not hard to show that an inscribed square takes up more than $1/2$ the area of a circle, and an inscribed regular octagon takes up more than $3/4$ of the area of the circle. Indeed, as the

ancient Greeks realized (Euclid's *Elements* XII 2), one can use what we call 'mathematical induction' to show that an inscribed regular 2^n -gon takes up more than $1 - 1/2^{n-1}$ of the area of a circle.

If we inscribe a regular 2^n -gon in a circle, its longest diagonals are diameters of that circle. The Greeks knew that the area of a regular 2^n -gon is proportional to the square on its longest diagonal – a result which follows from the fact that the area of a triangle with given angles is proportional to the square of its longest side – and from this it follows that, insofar as a circle is like a regular 2^n -gon, its area is proportional to the square on its diameter. Indeed, somewhat later than Antiphon, Eudoxus (408 – ca. 355 BC) gave a rigorous proof that the area of a circle is, in fact, proportional to the square on its diameter.

Antiphon boldly claimed that a circle simply *is* a regular polygon (with a large number of sides). In making this claim, Antiphon entered a lively discussion, started by Zeno (450 BC) and others, about whether space is continuous or discrete. If space is discrete, then there is some minimum area e . If n is so large that $1/2^{n-1}$ of the area of the circle is less than e , then an inscribed regular 2^n -gon, in taking up more than $1 - 1/2^{n-1}$ of the area, actually takes up *all* the area.

Another early Athenian mathematician who worked on the geometry of the circle was Hippocrates, who came from the Greek island of Chios, near present-day Turkey. (He is not to be confused with the physician, famous for his oath, who came from Cos.) Hippocrates, it is said, had been swindled in business and came to Athens about 430 BC to recover his property through legal action. The case dragged on, and Hippocrates used the time to study philosophy and supported himself by teaching geometry.

Hippocrates was responsible for much of the material in Books III and IV of Euclid's *Elements*. He called the square of a quantity 'dynamis', hence our 'power'. He pioneered the custom of reducing one theorem to another and may have been one of the first to use the method of *reductio ad absurdum* in mathematics.

He was also the first to find the precise area of a region bounded by curves, as we shall now see. Construct semicircles on three sides of a right triangle. By the converse of the theorem of Thales, the semicircle on the hypotenuse passes through the vertex at the right angle. The semicircles on the other two sides of the right triangle are supposed to lie outside the triangle. (See Figure 12.1.) The areas included in the two smaller semicircles, but not in the semicircle on the hypotenuse, are called *lunes* (after the crescent moon).

Hippocrates argued as follows. If the vertices of the triangle are A , B and C , with the right angle at C , then $AC^2 + CB^2 = AB^2$ (by the theorem of Pythagoras). Since the area of a circle is proportional to the square on its diameter, the area of a semicircle is likewise proportional to the square on

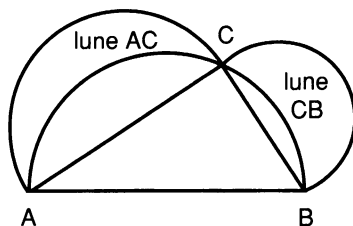


FIGURE 12.1. Lunes of Hippocrates

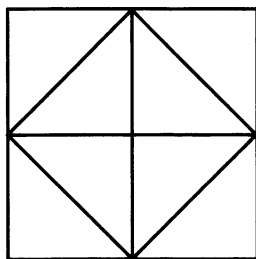


FIGURE 12.2. Doubling the square

the diameter. Therefore the sum of the areas of the semicircles on AC and CB equals the area of the semicircle on AB . Subtracting the areas where the semicircles overlap, we may conclude that the sum of the areas of the lunes equals the area of the right triangle. That is, since semicircle AC plus semicircle CB equals semicircle AB , it follows that

$$\text{lune } AC + \text{lune } CB = \text{triangle } ABC = \frac{1}{2}BC \times AC.$$

Hippocrates also contributed to the problem of ‘doubling the cube’. This was the problem of determining the length x such that $x^3 = 2$, preferably using only the geometry of straight lines and circles.

Legend has it that, during the plague of 430 BC, the Athenians consulted the oracle of Delos for help. The oracle replied that they should double the altar of Apollo, a marble cube, in size. When the plague refused to abate, the oracle explained that the Athenians had doubled the edges of the cube, not its volume. Although the Athenians did not succeed with this task, at least not according to the methods acceptable to Plato, the plague seems to have stopped anyway.

Hippocrates noted that one could double the volume of a cube, with edge one unit in length, if one could find quantities x and y such that $1/x = x/y = y/2$; for then x^3 would equal 2. However, he did not succeed in constructing these quantities in a way that satisfied Plato (who wanted to use only straightedge and compass).

Socrates (469–399 BC) was the mentor of Plato. As Plato portrays him in the dialogue *Meno*, Socrates claimed that all knowledge is recollection. In an argument with Meno on the nature of virtue, Socrates bet him that he could make his slave ‘remember’ a geometric construction and its proof. Asked to double the unit square, the slave, who was completely ignorant of geometry, first offered to double its side, but was soon led to admit his error. Then Socrates got him to look at the figure of a square with the midpoints of its four sides all joined to each other (Figure 12.2) and soon persuaded him to ‘remember’ that it is the square on the diagonal of the inner square which has double its area.

As a young man, Plato (429–349 BC) was a disciple of Socrates. After the latter’s death, Plato travelled to Africa, where he visited Heliopolis, now a suburb of Cairo, and Cyrene in Lybia. There he studied with Theodorus, who had proved the irrationality of the square roots of the nonsquare integers less than 18. Plato also went to Italy and became acquainted with Archytas (428–347 BC), the head of the Pythagorean school. Archytas also had ‘doubled the cube’, but he did so by going beyond the geometry of straight lines and circles.

Plato returned to Athens in about 380 BC and founded the famous Academy. At the entrance of this school was the inscription: *let no one ignorant of geometry enter here*.

The importance of Plato in the history of mathematics is due not so much to any mathematical contribution of his own as to the influence he exerted on others. It was he who insisted that a ‘proper’ solution involve no curves other than the circle (*Timaeus* 34a). It was he who emphasized the importance of clear definitions and postulates. Finally, Plato strongly encouraged people to study mathematics because he believed that this study would help them become wise and therefore virtuous. The five Platonic solids, or regular polyhedra, were not discovered by Plato, but he discussed them in the *Timaeus*.

Plato had a brilliant student, Theaetetus, who died in battle in 369 BC, and to whom he dedicated a dialogue. It was Theaetetus who showed that the square root of a natural number is irrational if and only if the natural number is not a square (*Theaetetus* 147c–148b). Theaetetus also studied the regular polyhedra, and worked on the theory of proportion. According to van der Waerden [1985], Theaetetus was responsible for Books X and XIII of Euclid’s *Elements*.

The most important Athenian mathematician at this time was Eudoxus of Cnidus, another small Greek island near modern Turkey. Eudoxus lived from 408 to 355 BC, and distinguished himself in astronomy, medicine, geography and philosophy – as well as mathematics. Like Plato, he studied astronomy in Heliopolis, and mathematics with Archytas in Tarentum (in what is now southern Italy). As a young man, Eudoxus studied in Plato’s

Academy, commuting on foot from Piraeus, the harbour district. Later he engaged in a philosophical controversy with Plato; it seems that Eudoxus anticipated the Epicurean position that humans strive to maximize pleasure minus pain.

In mathematics, Eudoxus was responsible for Books V and XII of Euclid's *Elements*. Book V deals with the theory of proportion. Today, might define the proportion ' a is to b as c is to d ' (written $a : b :: c : d$) as an equation $a/b = c/d$, and we would say that the proportion held just in case $ad = bc$. However, this presupposes our theory of the real number field. It presupposes that we already have some way of understanding what it is to multiply two irrational numbers. Eudoxus was starting from scratch. He could not use multiplication to define proportion because it was in terms of proportion that he defined multiplication. Eudoxus used his theory of proportion to prove the basic laws of multiplication, such as commutativity and associativity. The definition on which he based his development of the number system was the following:

$a : b :: c : d$ if and only if, for all positive integers p and q ,

$$pa > qb \text{ if and only if } pc > qd,$$

and likewise with $>$ replaced by $<$.

If we take a/b and c/d to be positive real numbers, this statement asserts:

- the set of rationals above a/b = the set of rationals above c/d
- the set of rationals below a/b = the set of rationals below c/d ,

thus anticipating the modern definition of real numbers due to Dedekind.

Actually, Eudoxus assumed that a, b, c and d are geometric quantities. For example, a and b could be arcs of circles and c and d could be angles. This is another reason why he did not write $ad = bc$; for how do you multiply an arc by an angle?

One particular ratio Eudoxus was interested in arose from the following problem; to divide a segment AB by a point H so that $AB/AH = AH/HB$, that is, the whole is to the larger part as the larger part is to the smaller. Taking $AH = x$ and $HB = 1$, we obtain the quadratic equation $x^2 - x - 1 = 0$, so that, upon discarding the negative solution, we find $x = (1 + \sqrt{5})/2$. This, or sometimes its reciprocal $(-1 + \sqrt{5})/2$, is known as the *golden section*.

Eudoxus also gave a proof that the area of a circle is proportional to the square on its diameter (Euclid's *Elements* XII 2), by inscribing regular polygons with 2^n sides in both circles and taking n sufficiently large.

Finally, we should mention Menaechmus (350 BC), another student of Plato's. Menaechmus discovered the conics — the ellipse, hyperbola and parabola — and used them to 'double the cube'. Using modern analytic

geometry, we can express his solution to the ‘Delian problem’ in the following simple fashion: the parabolas $y = \frac{1}{2}x^2$ and $x = y^2$ intersect at a point whose y coordinate is the cube root of 2. In his ‘doubling of the cube’, Menaechmus did not stick to straight lines and circles. Indeed, as already mentioned, in 1837, Pierre Wantzel showed that it is not possible to obtain a segment of length equal to the cube of root of 2, using only the geometry of straight lines and circles. We shall give a proof of this in Chapter 14.

Exercises

1. Prove that the inscribed regular 2^n -gon takes up more than $1 - 1/2^{n-1}$ of the area of the circle.
2. If the diameter of a circle is d , prove that the area of the inscribed regular 2^n -gon is $2^{n-3}d^2 \sin(\pi/2^{n-1})$.
3. Prove the theorem of Theaetetus that a natural number has an irrational square root if and only if it is not a perfect square.
4. Using the definition of proportion given by Eudoxus, show that $a : b :: c : d$ if and only if $d : c :: b : a$.

Plato and Aristotle on Mathematics

Plato (427–347 BC) believed that the objects in the universe fall into two very different classes, the material and the immaterial. A chair or an ox belongs to the class of material things. A soul or a number belongs to the class of immaterial things. The drawing of a square belongs to the material realm but the square itself belongs to the immaterial realm. Plato says of the students of geometry that they

make use of the visible forms and talk about them, though they are not thinking of them but of those things of which they are a likeness, pursuing their inquiry for the sake of the square as such and the diagonal as such, and not for the sake of the image of it which they draw (*Republic* 510d).

For Plato, the class of material things is characterized by change, uncertainty, ignorance and imperfection. The drawing of a square can be erased and it is doubtful whether its angles are each exactly 90° or whether its sides are perfectly straight.

On the other hand, the class of immaterial things is characterized by their constancy and perfection and by our certain knowledge of them. The square ‘as such’ has sides which remain perfectly straight forever. Its properties can be deduced with infallible rigour. We can know with absolute certainty that its diagonals are equal.

Scientists understand change and motion in the universe in terms of unchanging formulas or laws. Plato had a similar outlook (*Timaeus* 52-58). Moreover, he stressed that the formulas or laws have an existence of their own, independent of the material universe.

According to Plato, mathematical objects are not the only immaterial objects. Other immaterial objects are God, goodness, courage and the human soul (*Republic* 380d-383c). However, the best way to begin to know the immaterial realm is to do mathematics. One is to study number theory 'for facilitating the conversion of the soul itself from the world of generation to essence and truth' (*Republic* 525c). One is to study geometry 'to facilitate the apprehension of the idea of good' (*Republic* 525e).

Plato believes that the truths of mathematics are absolute, necessary truths. He believes that, in studying them, we shall be in a better position to know the absolute, necessary truths about what is good and right, and thus be in a better position to become good ourselves.

Platonism as a philosophy of mathematics is the view that at least the most basic mathematical objects (e.g., real numbers, Euclidean squares) actually exist, independently of the human mind which conceives them. Their properties are discovered, not created.

Aristotle (384–322 BC) was a student of Plato, but he disagreed with him about the nature of mathematics. In Book XIII of the *Metaphysics*, Aristotle asserts that

conclusions contrary alike to truth and to the usual views follow,
if one is to suppose the objects of mathematics to exist thus as
separate entities (*Metaphysics* 1077a).

For Aristotle, a word like 'two' is not a noun designating an abstract object but rather an adjective describing a concrete object (the *two* yard ladder, a *two* year period).

Whereas Platonism is quite compatible with the view that there are actually infinite lines and sets with an infinite numbers of elements, Aristotle is a staunch finitist. He would have rejected Cantor's 'aleph-null' (*Metaphysics* 1084a); he would have rejected infinitesimals (*Physics* 266b); he did reject infinite sets and infinite magnitudes (*Physics* III). For Aristotle, the geometer can have as much as he needs of an infinite line but he cannot have the whole line in its infinite totality.

Under the influence of Plato, Aristotle formulated a principle according to which every (mathematical) statement is either true or false, but he had his doubts when it came to applying this principle to the everyday temporal world. He wondered whether a statement like 'there will be a sea battle tomorrow' is either true or false (*On Interpretation* 9). How can it be true if the battle may not occur? How can it be false if the fight is a real possibility?

Like the view of the 20th century 'intuitionists', Aristotle's view is human-centered. The reality of numbers has to do, not with some alien heaven, but with the way we describe our surroundings. The infinite must be rejected because we humans work in a finite way. The truth about certain

propositions may be left in abeyance if it is, for the moment, inaccessible to human beings.

Aristotle also had something to contribute to the old problem that had divided the Ionic philosophers from those based in Italy: whether the universe is made up of substances or atoms, whether things should be measured or counted. He pointed out that when we talk about a loaf of bread or a glass of wine, bread and wine are measured, but loaves and glasses are counted. He asserted that we measure *matter* by counting its *forms*.

Exercise

How might Aristotle answer Zeno's arguments against motion?

14

Constructions with Ruler and Compass

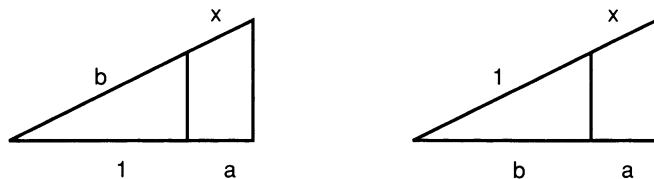
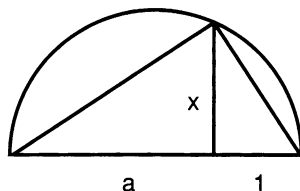
Ancient Greek mathematicians were haunted by three problems:

- I** doubling the cube, that is, finding the cube root of 2;
- II** trisecting any given angle, say an angle of 60°
(of course some angles are easily trisected, for example, one of 90°);
- III** squaring the circle,
that is, constructing a square equal in area to that of a given circle.

Assorted solutions to these problems were proposed at various times, but these did not conform to the rules of the game, presumably laid down in Plato's Academy, that *only constructions with ruler and compass* be admitted. (Actually, we only know that Pappus attributed these rules to Plato more than 600 years later.) Moreover, the ruler could only be used for joining two points and the compass could only be used for drawing a circle with a given point as center and a given segment as radius. The reader will have no difficulty in carrying out the following constructions with ruler and compass:

- (a) to bisect a given angle;
- (b) to find the right bisector of a given segment;
- (c) to draw a line through a given point parallel to a given line;
- (d) to construct an equilateral triangle.

If we adopt a given segment as our unit of length, we can represent any positive real number by a segment, actually by the ratio of this segment to

FIGURE 14.1. Finding ab and a/b FIGURE 14.2. Constructing root of a

the unit segment. With the help of ruler and compass, the Greeks were able to perform the following arithmetical operations on positive real numbers: adding, subtracting (the smaller from the larger), multiplying, dividing and extracting square roots. The first four of these are called *rational* operations.

Indeed, for addition and subtraction this is obvious. To find $x = ab$ one considers the proportion

$$x : b = a : 1$$

and to find $x = a/b$ one considers the proportion

$$x : 1 = a : b.$$

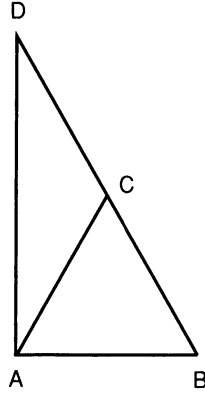
In both cases, the problem is that of finding the fourth proportional to three given lengths, which can easily be done, following Thales, with the aid of similar triangles, using only ruler and compass constructions. See Figure 14.1.

To find $x = \sqrt{a}$, one looks at the proportion

$$1 : x = x : a.$$

Here the problem is that of finding the *mean* proportional of two given lengths, which the Greeks solved ingeniously by ruler and compass constructions, as illustrated by Figure 14.2, which exhibits the semicircle on a segment of length $a + 1$.

To attack problem II, the trisection problem, we would nowadays use trigonometry. First note that we can construct an angle if and only if we can construct its cosine. If $\theta = 60^\circ/3 = 20^\circ$, then $\cos 3\theta = \cos 60^\circ = 1/2$, as is seen from Figure 14.3.

FIGURE 14.3. The cosine of 60°

In Figure 14.3, ABC is an equilateral triangle and $CD = AC$. It easily follows that the angle at B is 60° , the angle at D is 30° and the angle BAD is 90° , hence $\cos 60^\circ = AB/BD = 1/2$.

Now

$$\begin{aligned}
 \cos 3\theta &= \cos(2\theta + \theta) = \cos 2\theta \cos \theta - \sin 2\theta \sin \theta \\
 &= (\cos^2 \theta - \sin^2 \theta) \cos \theta - 2 \sin \theta \cos \theta \sin \theta \\
 &= \cos^3 \theta - 3 \sin^2 \theta \cos \theta = \cos^3 \theta - 3(1 - \cos^2 \theta) \cos \theta \\
 &= 4 \cos^3 \theta - 3 \cos \theta.
 \end{aligned}$$

Thus we want to solve the equation $8 \cos^3 \theta - 6 \cos \theta = 1$. Putting $2 \cos \theta = u$, we obtain the cubic equation

$$u^3 - 3u - 1 = 0.$$

The question is therefore whether a solution of this cubic equation can be expressed in terms of rational operations and square roots. We shall return to this problem in the next chapter. First let us look at a problem which the Greeks were able to solve by their methods.

One of the highlights in Euclid's *Elements* is the construction of a regular pentagon, equivalently, that of an angle of $360^\circ/5 = 72^\circ$. Today we would attack this problem too with the help of trigonometry; for an elegant argument we may even invoke complex numbers. Let $\theta = 72^\circ$, then $5\theta = 360^\circ$, hence, by de Moivre's Theorem,

$$\begin{aligned}
 (\cos \theta + i \sin \theta)^5 &= \cos 5\theta + i \sin 5\theta \\
 &= \cos 360^\circ + i \sin 360^\circ \\
 &= 1 + i0 \\
 &= 1.
 \end{aligned}$$

Thus, we wish to solve the equation $z^5 = 1$, that is,

$$(z - 1)(z^4 + z^3 + z^2 + z + 1) = 0,$$

where $z = \cos \theta + i \sin \theta$. Note that $z^{-1} = \cos \theta - i \sin \theta$, so that

$$2 \cos \theta = z + z^{-1} = u$$

say. Clearly, $z = 1$ is not a satisfactory solution, so our problem reduces to solving

$$z^4 + z^3 + z^2 + z + 1 = 0.$$

Dividing by z^2 , we write this as

$$z^2 + z + 1 + z^{-1} + z^{-2} = 0.$$

Now $z + z^{-1} = u$, hence $z^2 + 2 + z^{-2} = u^2$, so that $z^2 + z^{-2} = u - 2$. The equation to be solved then becomes

$$u^2 + u - 1 = 0.$$

Discarding the negative solution, we find

$$u = \frac{-1 + \sqrt{5}}{2},$$

a number which we recall from Chapter 12 as the golden section.

Of course this is not how the Greeks attacked the problem, as they did not know about complex numbers, and, at the time of Euclid, had not yet invented trigonometry. Nonetheless, our analysis shows that they could construct the angle θ , since $2 \cos \theta = (-1 + \sqrt{5})/2$ involves only rational operations and square roots. Indeed, when Euclid constructs a regular pentagon in Book IV, Proposition 11, by what looks to us like a rather complicated method, he makes use (via Proposition 10) of the earlier construction of the golden section in Book II, Proposition 11 (taken up once more in Book VI, Proposition 30). The Greeks could of course construct squares and regular hexagons, but one problem they did leave open was the following:

IV constructing a regular heptagon, that is, a seven sided figure.

By the same method we just employed for constructing a regular pentagon, we can see that this problem reduces to the following cubic equation in $u = 2 \cos(360^\circ/7)$:

$$u^3 + u^2 - 2u - 1 = 0.$$

Although the ancient Greeks did not know this, the only arithmetical operations that can be carried out with ruler and compass constructions are rational operations and square root extractions and, of course, combinations of such. To understand why this is so, we have to make use of analytic geometry, which was only developed in the seventeenth century by René Descartes.

His pioneering idea was to represent every point in the plane by a pair of real numbers (x, y) and to observe, conversely, that every such pair represents a point. Unlike the Greeks, we need not confine x and y to be positive: if we use the modern rectangular coordinate system, x is negative in the second and third quadrant, y is negative in the third and fourth quadrant. We can now say that a straight line consists of all points (x, y) satisfying an equation of the form

$$ax + by + c = 0,$$

where a, b and c are given real numbers, and a circle consists of all points (x, y) satisfying an equation of the form

$$x^2 + y^2 + dx + ey + f = 0.$$

Now what happens when we perform the following operations:

1. join two given points,
2. draw a circle with given center and radius,
3. intersect two straight lines,
4. intersect a circle and a straight line,
5. intersect two circles?

(1) Suppose the given points are (x_1, y_1) and (x_2, y_2) . Then we easily see that the straight line has the equation

$$(y_1 - y_2)x + (x_2 - x_1)y + (x_1y_2 - x_2y_1) = 0,$$

in other words, an equation of the form

$$ax + by + c = 0,$$

where a, b and c are expressed in terms of the given quantities x_1, y_1, x_2, y_2 by means of the operations of addition, subtraction and multiplication.

(2) Suppose the center is (α, β) and the radius is ρ , then the equation of the circle is

$$(x - \alpha)^2 + (y - \beta)^2 = \rho^2,$$

that is,

$$x^2 + y^2 - 2\alpha x - 2\beta y + \alpha^2 + \beta^2 - \rho^2 = 0.$$

Thus the equation of the circle is of the form

$$x^2 + y^2 + dx + ey + f = 0,$$

where d, e and f are again expressed in terms of the given quantities α, β and ρ by means of addition, subtraction and multiplication.

(3) To find the intersection of two given straight lines, we must solve the pair of equations:

$$ax + by + c = 0,$$

$$a'x + b'y + c' = 0.$$

We obtain the solution

$$x = -\frac{cb' - c'b}{ab' - a'b}, \quad y = -\frac{ac' - a'c}{ab' - a'b}.$$

(It is assumed that $ab' - a'b \neq 0$, otherwise the two lines are parallel or even coincide.) Again we see that the new quantities x and y are obtained from the given quantities a, b, c, a', b' and c' by means of the rational operations, including division.

(4) To find the intersection of a circle and a straight line, we must solve the pair of equations:

$$ax + by + c = 0,$$

$$x^2 + y^2 + dx + ey + f = 0.$$

Assuming, for example, that $b \neq 0$, we get

$$y = -\frac{a}{b}x - \frac{c}{b}$$

from the first equation. When we substitute this into the second equation, we obtain a quadratic equation:

$$Ax^2 + Bx + C = 0,$$

where A, B and C are expressed by means of the rational operations in terms of the given quantities a, b, c, d, e and f . In particular, $A = 1 + a^2/b^2 > 0$. Finally, solving for x , we obtain:

$$x = \frac{-B \pm \sqrt{B^2 - 4AC}}{2A}.$$

Here x is expressed not only by means of the rational operations, but also requires a square root. Note that, if $B^2 - 4AC = 0$, the line is tangent to the circle and, if $B^2 - 4AC < 0$, it does not meet the circle at all.

(5) To find the intersection of two circles, we must solve the pair of equations:

$$\begin{aligned}x^2 + y^2 + dx + ey + f &= 0, \\x^2 + y^2 + d'x + e'y + f' &= 0.\end{aligned}$$

By subtracting, we may here replace the second equation by the linear equation

$$(d - d')x + (e - e')y + f - f' = 0,$$

so the situation is the same as that already treated under (4) above. The last equation represents the straight line passing through the two points of intersection of the given circles, or, if the given circles merely touch, their common tangent. Of course, it may also happen that the two circles have no common points at all.

We had seen earlier that geometric constructions with ruler and compass allow us to carry out the rational operations and to extract square roots. We have now shown the converse: combinations of rational operations and square roots are the only arithmetical operations which can be carried out in this way. That is, when we perform ruler and compass constructions (1) to (5), we only get lengths that can be expressed in terms of the operations $+$, $-$, \times , $/$, and $\sqrt{}$. To solve problems **I** to **IV** by ruler and compass constructions is thus equivalent to expressing the real numbers $\sqrt[3]{2}$, $\cos 20^\circ$, $\sqrt{\pi}$ and $\cos(360^\circ/7)$ by rational operations and square roots.

Exercises

1. Show how to carry out the constructions (a) to (d) in the text, using ruler and compass only.
2. Carry out a ruler and compass construction of the golden section $(-1 + \sqrt{5})/2$.
3. To construct a regular heptagon one has to find the angle $\theta = 360^\circ/7$. Show that $u = 2 \cos \theta$ satisfies the cubic equation

$$u^3 + u^2 - 2u - 1 = 0.$$

The Impossibility of Solving the Classical Problems

The ancient Greeks were unable to solve problems **I** to **IV** using ruler and compass constructions, for a good reason: it cannot be done. Concerning problem **III**, this was shown only in 1882 by C.L.F. Lindemann (1852–1939), who proved that π is not an algebraic number, which implies, in particular, that $\sqrt{\pi}$ cannot be constructed by rational operations and square roots. His method was based on an earlier proof by Hermite, who had shown that $e = \lim_{n \rightarrow \infty} (1 + \frac{1}{n})^n$ is not algebraic.

Problems **I**, **II** and **IV** have one thing in common: they can all be expressed by cubic equations, namely,

$$u^3 - 2 = 0, \quad u^3 - 3u - 1 = 0, \quad u^3 + u^2 - 2u - 1 = 0.$$

(For the last see Exercise 3 of Chapter 14.) First let us make sure that they have no rational solutions. (For the first equation we already know this.)

Lemma 15.1. *If $a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0 = 0$ is a polynomial equation with integer coefficients, then any rational solution is of the form p/q , with p a factor of a_0 and q a factor of a_n . In particular, when $a_n = 1$, any rational solution will be an integer and a factor of a_0 .*

Proof: Let p/q be a rational solution, where $q \neq 0$ and $\gcd(p, q) = 1$. Putting $x = p/q$ in the equation and multiplying by q^n , we obtain

$$a_n p^n + a_{n-1} p^{n-1} q + \cdots + a_0 q^n = 0.$$

Since p divides all the terms except possibly the last, it must divide the last term also. Since $\gcd(p, q) = 1$, it follows from the unique factorization into primes that p divides a_0 . Similarly, q divides a_n .

It follows from the lemma that the only possible rational solutions of the equation $u^3 - 2 = 0$ are $u = \pm 1$ or $u = \pm 2$, all four of which are quickly seen not to be solutions. Similarly, the only possible solutions of the second and third equations are $u = \pm 1$, which are also seen not to work. We may conclude that solutions to these three equations cannot be expressed by means of rational operations alone, but there remains the possibility that they can be expressed with the help of square roots.

Over the years, many people have tried their hand at problems **I** and **II**. For example, the much envied Casanova worked on doubling the cube and, even today, there are still many determined angle trisectors. However, in 1837, Pierre Wantzel (1814–1848) showed that none of $\sqrt[3]{2}$, $2 \cos 20^\circ$, and $2 \cos(370^\circ/7)$ can be expressed in terms of rational operations and square roots and, therefore, that problems **I**, **II** and **IV** cannot be solved by ruler and compass constructions.

To give a simple exposition of why this is so, we shall introduce a more modern concept, that of a field. For our purposes, a *field* is a set of numbers, real or complex, which contains the number 1 and which is closed under the rational operations. (There are other fields, but we shall not need them here.) In particular, the rationals form a field \mathbf{Q} and so does $\mathbf{Q}[\sqrt{2}]$, the set of all numbers of the form $a + b\sqrt{2}$, where a and b are rational. More generally, if F is any field, then so is $F[\sqrt{c}]$, where c is a given element of F , by which we understand the set of all numbers of the form $a + b\sqrt{c}$ with $a, b \in F$.

It is clear that $F[\sqrt{c}]$ is closed under addition, subtraction and multiplication. To show that it is also closed under division, we assume that \sqrt{c} is not in F , otherwise there would be nothing to prove, and that $a + b\sqrt{c} \neq 0$. We calculate

$$\frac{1}{a + b\sqrt{c}} = \frac{1}{a + b\sqrt{c}} \times \frac{a - b\sqrt{c}}{a - b\sqrt{c}} = \frac{a - b\sqrt{c}}{a^2 - b^2c} = \frac{a}{a^2 - b^2c} + \frac{-b}{a^2 - b^2c}\sqrt{c},$$

which is again of the form $a' + b'\sqrt{c}$ with $a', b' \in F$. A small argument is necessary to check that $a^2 - b^2c \neq 0$.

We can now say that a real number u is constructible with ruler and compass, equivalently, expressible by rational operations and square roots, if and only if there exists a sequence of fields

$$\mathbf{Q} = F_0 \subseteq F_1 \subseteq \cdots \subseteq F_n$$

such that $F_{k+1} = F_k[\sqrt{c_k}]$ with $c_k \in F_k$ and $u \in F_n$.

Proposition 15.2. *Suppose $f(x) = x^3 + a_2x^2 + a_1x + a_0$ is a cubic polynomial with coefficients in a field F . Suppose further that the equation $f(x) = 0$ has a solution in $F[\sqrt{c}]$ with $c \in F$. Then it already has a solution in F .*

Proof: Let $x_1 = a + b\sqrt{c}$ be the given solution with $a, b \in F$. Then

$(x_1 - a)^2 = b^2c$, hence $x_1^2 + px_1 + q = 0$ with $p, q \in F$. Dividing $f(x)$ by the polynomial $x^2 + px + q$, we obtain

$$f(x) = (x^2 + px + q)(x + d) + (ex + f),$$

where the quotient is $x + d$ with $d \in F$ and the remainder is $ex + f$ with $e, f \in F$. Since $f(x_1) = 0$ and $x_1^2 + px_1 + q = 0$, we deduce that $ex_1 + f = 0$. If $e \neq 0$, then $x_1 = -f/e \in F$ and we need look no further. If $e = 0$, then also $f = 0$, hence $x + d$ is a factor of $f(x)$. But then $f(-d) = 0$ and so $x_2 = -d \in F$ is the required solution.

Corollary 15.3. *If a number expressible by rational operations and square roots satisfies a cubic equation with rational coefficients, then this equation must have a rational solution.*

Proof: Suppose the cubic equation $f(x) = 0$ has no rational solution. Then, by Proposition 15.2 with $F = \mathbf{Q}$, it has no solution in $\mathbf{Q}[\sqrt{c_1}]$ with $c_1 \in \mathbf{Q}$. Again, by the Proposition with $F = \mathbf{Q}[\sqrt{c_1}]$, it has no solution in $\mathbf{Q}[\sqrt{c_1}][\sqrt{c_2}]$ with $c_2 \in \mathbf{Q}[\sqrt{c_1}]$. Continuing in this way, we see that it has no solution in $\mathbf{Q}[\sqrt{c_1}] \cdots [\sqrt{c_n}]$ for any n , where $c_k \in \mathbf{Q}[\sqrt{c_1}] \cdots [\sqrt{c_{k-1}}]$ for $1 < k \leq n$. Thus it has no solution expressible by rational operations and square roots.

Since the cubic equations

$$u^3 - 2 = 0, \quad u^3 - 3u - 1 = 0, \quad u^3 + u^2 - 2u - 1 = 0$$

have no rational solutions, as we verified earlier, we can now infer from Corollary 15.3 that they have no solutions expressible in terms of rational operations and square roots. In view of Chapter 14, we may therefore conclude that problems **I**, **II** and **IV** cannot be solved using only ruler and compass constructions. In summary:

Theorem 15.4. *It is impossible to double a cube, to trisect an arbitrary angle or to draw a regular heptagon by ruler and compass constructions.*

We have seen that the Greeks were able to draw regular polygons with 3 or 5 sides, but not with 7 sides. The question arises, for which primes p is it possible to construct a regular p -gon using ruler and compass only? Carl Friedrich Gauss showed that this is possible whenever p is a prime of the form $2^n + 1$ and Wantzel proved the converse. Gauss was so pleased with his discovery that he wanted a regular 17-gon inscribed on his tombstone. His request was not carried out, but a regular 17-gon was inscribed on a monument to Gauss in Braunschweig, Germany.

Odd prime numbers of the form $2^n + 1$ are called ‘Fermat primes’, after Pierre de Fermat (1601–1665). It is easy to prove that $2^n + 1$ cannot be prime unless n has the form 2^k .

To see this, one makes use of the identity

$$x^b + 1 = (x + 1)(x^{b-1} - x^{b-2} + \cdots + 1)$$

for odd b . If $n = ab$ and b is odd, put $x = 2^a$ and infer that $2^a + 1$ is a factor of $2^n + 1$. Unless $b = 1$, it follows that $2^a + 1$ is an odd proper divisor of $2^n + 1$ different from 1, so $2^n + 1$ cannot be prime. Thus, if $2^n + 1$ is prime, n cannot have any odd proper divisors other than 1, hence n must be a power of 2.

Fermat was under the impression that $2^{2^k} + 1$ is prime for all natural numbers k . Euler later found that $2^{2^5} + 1$, a ten digit number, is divisible by 641. An easy, though tricky, way of seeing this is as follows:

$$\begin{aligned} 2^{2^5} &= 16 \times 2^{28} = (641 - 5^4)2^{28} \\ &= 641m - (5 \times 2^7)^4 \\ &= 641m - (641 - 1)^4 \\ &= 641m - (641n + 1) \\ &= 641(m - n) - 1, \end{aligned}$$

where m and n are integers, hence $2^{2^5} + 1$ is a multiple of 641.

At the moment, the only known Fermat primes are 3, 5, 17, 257, and 65,537, corresponding to $k = 0, 1, 2, 3$ and 4, respectively. Not surprisingly, two people in the 19th century tried to break records by actually constructing regular polygons with 257 and 65,537 sides.

For $k = 5, 6, \dots, 22$, it is known that $2^{2^k} + 1$ is composite.

It follows from the work of Gauss and Wantzel that a regular polygon with m sides can be constructed by ruler and compass if and only if

$$m = 2^k p_1 \cdots p_l,$$

where $k \geq 0$ and p_1, \dots, p_l are distinct Fermat primes.

Exercises

1. If $a, b \in F$ but $\sqrt{c} \notin F$, show that $a + b\sqrt{c}$ is either 0 or possesses an inverse in $F[\sqrt{c}]$.
2. If $d \in F[\sqrt{c}]$, show that any element of $F[\sqrt{c}][\sqrt{d}]$ satisfies an equation of degree 4 with coefficients in F .
3. Explain how the Greeks could construct regular polygons with 15 and 60 sides.
4. If n is a positive integer, show that an angle of n degrees can be constructed with ruler and compass if and only if n is a multiple of 3.

Euclid

The city of Alexandria, on the mediterranean coast of Egypt, was founded by Alexander the Great in 332 B.C., who brought Greeks, Egyptians and Jews to settle there. One of his generals, Ptolemy I, made Alexandria the capital of his kingdom and founded a dynasty consisting of a long line of rulers, also named 'Ptolemy' and ending with the reign of the famous queen Cleopatra, who picked the wrong side in a Roman civil war.

Ptolemy established a university in Alexandria, called the 'Museum', which was soon to acquire a library holding more than 600,000 papyrus scrolls. For well over 600 years, Alexandria was to be the mathematical and scientific center of the world, with only some schools of philosophy surviving in Athens, although, after the extinction of the Ptolemaic line with Cleopatra, Alexandria was ruled by Rome. It was ultimately conquered by the Arabs in 641 AD.

The first chair of mathematics at the Museum was occupied by Euclid (330 to 275 BC), said to have been a student of a student of Plato. Apart from a couple of anecdotes, we know little about his life, and some ancient authors even thought he was a committee, like the 20th century Nicolas Bourbaki. According to one anecdote, Euclid told the impatient king that 'there is no royal road to learning'. According to another, he gave a small coin to a student who demanded to know the practical value of the lectures he had been attending.

Euclid wrote a number of books, on optics, music, astronomy etc., but his fame rests on the *Elements*, a collection of 13 so-called books (which we would now call chapters), which presented the foundations of all the mathematics known in his day. Nothing like this was to be published again, until

the middle of the 20th century, when Nicolas Bourbaki issued a collection of books that purported to cover the elements of all the mathematics we study now.

None of the theorems contained in the 13 books can with certainty be ascribed to Euclid himself. It is believed that the Pythagoreans, including Archytas, were responsible for much of what appears in Books I, II, VI, VII, VIII, IX and XI and that Hippocrates was behind Books III and IV. For Books V and XII we are to thank Eudoxus, and Books X and XIII are said to be based on the work of Theaetetus.

However, the logical organization of the *Elements* is undoubtedly Euclid's contribution. Its success can be measured by the fact that, after more than 2,000 years, it was still used as a textbook in British schools. Moreover, throughout the ages, its structure was often imitated. Thomas Aquinas used a similar axiomatic presentation in his *Summa*, Newton's *Principia* is written in the style of the *Elements* and Spinoza's *Ethics* follows its logical arrangement. Undoubtedly the *Elements* has been the most influential scientific textbook in history.

Euclid's grandiose plan was to deduce all of mathematics from a small number of initial definitions and assumptions. The assumptions are subdivided into *axioms*, dealing with mathematics in general, and *postulates*, dealing with geometry in particular.

His treatment illustrated the ideal described by Aristotle at the beginning of his *Posterior Analytics*: sure knowledge is obtained by the rigorous deduction of the consequences of basic truths. To Euclid, these basic truths were either definitions or basic assumptions, largely assertions of unique existence. Let us take a closer look at his definitions, axioms and postulates.

The *Elements* begins with a list of 23 definitions, of which we will mention the first four:

1. A *point* is that which has no parts.
2. A *line* is length without width.
3. The extremities of a line are points.
4. A *straight line* is a line which lies evenly with the points on itself.

These statements are not definitions in the modern sense, though they make it clear that a point has no extension, that a line is not necessarily straight and that it is of finite length. Today we prefer to regard points and straight lines as undefined primitive concepts and leave the definition of curved lines to more advanced mathematics. The obscurity of Definition 4 may be due to the translation.

Euclid's axioms are intended to apply to all of mathematics, not just to geometry. A typical axiom asserts: 'If equals are added to equals, their sums are equal.' One cannot quarrel with this statement, though today

we might derive it from axioms of equality and the view of addition as an operation.

Euclid lists five postulates, which we shall now state and comment upon.

I. To draw a straight line from any point to any other point.

Presumably this means that there exists a unique straight line joining two distinct given points. Thus, a 'straight line' cannot be interpreted as referring to a great circle on a sphere, as there are many great circles joining two antipodes, e.g., the meridians passing through the two poles on the globe. The way to get around this objection is to *identify* antipodal points; one then obtains *elliptic geometry*, which also satisfies Postulate I.

II. To produce a finite straight line continuously in a straight line.

Here 'continuously' is usually interpreted to imply 'indefinitely', thus ruling out not only spherical, but also elliptic geometry.

III. To describe a circle with any center and any distance [as radius].

Like Postulate I, this is a construction, or unique existence statement, the word 'circle' having previously been defined, in Definition 15, as 'a plane figure contained by one line [i.e. curve] such that all the straight lines falling upon it from one point among those lying within the figure are equal to one another.' It would appear that by a circle Euclid means not just its circumference but also its interior.

IV. That all right angles are equal to one another.

The status of this assertion as a postulate is rather dubious, and it has been argued, already in antiquity, that it should be listed as an axiom instead.

V. That, if a straight line falling on two straight lines makes the interior angles on the same side less than two right angles, the two straight lines, if produced indefinitely, meet on that side on which are the angles less than two right angles.

This is the most famous of Euclid's postulates and it is to his credit that he recognized its significance. It will be discussed at length in Chapter 17.

It is on these definitions and assumptions that Euclid plans to erect his impressive edifice of logical deductions. Here is how he begins:

Proposition 1.

On a finite straight line to construct an equilateral triangle.

In his proof he considers a segment AB and constructs circles with centers A and B and radius AB . He then considers the point C in which the two circles intersect and goes on to show that $\triangle ABC$ is equilateral.

This proof falls short of modern standards of rigour. In general, two circles may meet in two points, touch at one point, or not meet at all. In the present situation, they do, in fact, meet in two points; but this does not follow from Euclid's explicit assumptions.

Book I concludes with proofs of the Theorem of Pythagoras and its converse. Euclid is careful to show that there is a square on the hypotenuse, before discussing its properties. This is interesting, in view of Legendre's later proof that the existence of such a square implies Euclid's Postulate V.

To prove the Theorem of Pythagoras, Euclid uses a theory of area. Nowadays we are tempted to define the area of a rectangle as 'length times width'. This presupposes a theory which explains what it means to multiply two irrationals. Euclid approached the question of area from a more elementary point of view. He began with the idea that two polygons have the same area if they first can be dissected into triangles which can be reassembled, as in a jigsaw puzzle, to form a polygon exactly like the second polygon. It is only in Book VI, after Euclid has presented Eudoxus's theory of irrationals, that the length times width formula is justified.

However, in Book II, Euclid gives geometric treatments of certain basic algebraic identities, such as $a(b+c) = ab+ac$, using the areas of rectangles to handle products. He also gives a proof of a statement equivalent to what we now call the Law of Cosines.

Book III discusses the basic properties of the circle. Euclid goes to great length to give rigorous proofs. For example, in spite of the fact that it is 'obvious from the diagram', Euclid offers a demonstration of the fact that the points on a chord of a circle lie in the interior of the circle. Euclid is not always successful in his attempt at rigour, but it is clear that he does understand the need for it.

Book IV gives constructions for various regular polygons. It culminates with a treatment of the regular 15-gon. This achievement remained unsurpassed until 1796, when Carl Friedrich Gauss (1777–1855) found a construction for the regular 17-gon.

In Book V, Euclid uses Eudoxus's definitions of proportion to deduce an arithmetic for line segments. The 'commutativity of multiplication' is the subject of Proposition 16.

In Book VI, Euclid uses the material of Book V to derive the basic properties of similar triangles. The book concludes with the theorem that the length of a circular arc is proportional to the angle it subtends at the center of the circle. In talking about 'arclength', Euclid is implicitly presupposing the 'completeness' of the plane.

Books VII to IX present some elementary theorems of number theory. Included are proofs for Euclid's Algorithm (VII 2), the unique factorization of square-free integers (IX 14), the infinitude of primes (IX 20), the formula for the sum of a geometric progression (IX 35), and the formula for even perfect numbers (IX 36).

Book X is occupied with what we might call ‘field extensions of degree 4 over rationals’. Euclid is interested in knowing when an expression like $\sqrt{7 + 2\sqrt{6}}$, which looks like it has ‘degree 4’ is actually equal to an expression like $1 + \sqrt{6}$, which involves only one ‘layer’ of square roots.

Book XI derives the basic theorems of solid geometry. A ‘cone’ is defined in terms of the revolution of a right triangle. A ‘cube’ is ‘a solid figure contained by six equal squares’. Proposition XI 21 says that ‘any solid angle is contained by plane angles [whose sum is] less than four right angles’. This proposition is used at the end of Book XIII to show that there are at most 5 regular polyhedra.

Book XII is the masterpiece of Eudoxus. Without the help of calculus, he manages to give a rigorous treatment of the volumes of the pyramid, cone and sphere.

Book XIII is the apex of the *Elements*. For each of the five regular polyhedra, Euclid derives the ratio of its side to the radius of the sphere in which it is inscribed. Although Euclid failed to give a complete theory of regular polygons – for example, the construction of the regular 17-gon is missing – he succeeded in giving a complete theory of regular polyhedra.

Euclid’s *Elements*, or watered down versions of it, was used for over 2,000 years in universities and schools to teach not only geometry but also rigorous thinking. Not long after World War II, a reaction against this program set in and educators decided that geometry was not the appropriate place for training in logic. Anyway, they argued, Euclid was not rigorous enough and Hilbert’s rigorous treatment (Chapter 17) was too cumbersome. So geometry was swept away in favour of ‘New Mathematics’. The French mathematician Dieudonné, one of the founding members of the Bourbaki group, suggested that linear algebra should replace what he contemptuously called ‘the theory of the triangle’.

Exercises

1. True or false? If two triangles have the same area, you can cut one of them up into little triangles, which can then be placed side by side to form a triangle congruent to the second triangle. Give a reference or a reason for your answer.
2. How did Euclid construct the regular 15-gon?

Non-Euclidean Geometry and Hilbert's Axioms

The parallel postulate V, Euclid's fifth postulate, seems less natural or convincing than the others. Ever since Euclid's time, people have felt that it ought to be deducible from Euclid's other postulates I to IV or from some logically equivalent set of axioms.

Noteworthy attempts to prove Euclid's fifth postulate were made by Proclus (410–85 AD), Saccheri (1667–1733), Thibault (1775–1822), and many others. We now know that these attempts were doomed to fail. Postulate V is independent of I to IV and one of Euclid's contributions to mathematics was his implicit recognition of this fact by presenting V as an axiom.

Given 'absolute geometry' (that is, the geometry based only on postulates I to IV) there are a number of important statements equivalent to the parallel postulate:

- Through a point not on a line there is exactly one line parallel to that line — Playfair.
- Every segment is a side of a square (with four right angles) — Legendre.
- Not every pair of similar triangles is congruent — Wallis.
- Every triangle has a circumcircle — Legendre.
- There is at least one triangle whose angle sum is 180° — Legendre.

(Heath notes that the first of these is due to Proclus.)

The search for a proof of the parallel postulate led to the discovery of many such equivalent statements, but each one was felt to be insufficiently 'self-evident' or 'basic' to count as a proper Euclidean axiom. What was really wanted was a deduction of postulate V from postulates I to IV alone.

Gauss may have been the first person to suspect the truth. In a letter to Franz Taurinus, written in 1824, Gauss says that he is sure that the parallel postulate cannot be proved.

Consider the alternative postulate:

(H) Through any point not on a line, there are at least two lines through that point and parallel to that line.

If we replace Euclid's parallel postulate by (H), we get the axioms of 'hyperbolic geometry'. It seems that Gauss believed hyperbolic geometry to be consistent.

The first person to publish results in hyperbolic geometry was the Russian N. I. Lobachevsky (1793–1856), of the University of Kasan, in 1829. In the same year, essentially the same results were discovered independently by the Hungarian J. Bolyai.

It was not until 1868 that it was proved that postulates I to IV do not imply postulate V. In that year, E. Beltrami (1835–1900) gave a Euclidean model for hyperbolic geometry. This showed that, if hyperbolic geometry contained any logical contradiction (for example, the assertion that both (H) and V are true), then that contradiction could be translated into a contradiction in Euclidean geometry. Since, presumably, there is no inconsistency in Euclidean geometry, there is none in hyperbolic geometry either.

In 1882, in the first article ever published in *Acta Mathematica*, Henri Poincaré (1854–1912) gave a sketch of a second Euclidean model for hyperbolic geometry. This model goes as follows. We interpret 'point' as 'point of the Cartesian plane in the interior of the unit circle $x^2 + y^2 = 1$ '. We interpret a 'line' to mean either a 'diameter of the unit circle (minus endpoints)', or else 'a circular arc in the interior of the unit circle and orthogonal to it'. (Two arcs are 'orthogonal' if they intersect at right angles.)

'Betweenness' is defined in the obvious way. Segment equality is defined as follows. Let AB be a 'segment', that is, part of a diameter or orthogonal arc. Let A^e be the endpoint of that diameter or orthogonal arc which is on A 's side of the diameter or arc. Let B^e be the other endpoint. Let

$$d(AB) = (AB^e/BB^e)(BA^e/AA^e),$$

where the segments on the right of this equation are ordinary Euclidean segments. Then two Poincaré segments AB and CD are 'equal' if and only if $d(AB) = d(CD)$.

The 'angle' between two Poincaré lines is the Euclidean angle between their tangents through the point where the lines meet. Angle equality is defined in the usual way.

Given the usual definition of the 'circle', it turns out, somewhat surprisingly, that 'circles' in the Poincaré sense are ordinary circles. It is just that their Poincaré centers are not where you would expect.

Using Steiner's geometry of 'inversion', one can prove that, under this interpretation, postulates I to IV are theorems in Euclidean geometry. However, postulate V, so interpreted, is not; through the center of the unit circle, for example, one can draw many 'lines' which do not meet a given orthogonal arc.

If V followed from I to IV, then V interpreted in the Poincaré sense, would both hold and not hold. We would have a statement about Euclidean lines and circles (related to the unit circle) which was both provable and disprovable relative to I to IV. Thus, assuming that I to IV are consistent, so are I to IV together with (H).

The reader can find the details of the 'proof by inversion' on pages 402 to 407 of volume 1 of Eves [1963]. Note that, although Eves uses logarithms in his proof, the reasoning is just the same if one drops them. (Eves uses logarithms because he wants the smallest 'distance' to be 0, not 1, and because of some results he aims to derive in a later section of his book.)

Nowadays people consider not only hyperbolic, but also 'elliptic geometry'. This was developed by Riemann in 1854, but is not to be confused with the more general 'Riemannian geometry', which we shall discuss below. In elliptic geometry, the straight lines are finite, and there are *no* parallels. A 'point' is like a pair of points on a sphere, and a 'line' is like a great circle on that sphere. Unfortunately, elliptic geometry does not satisfy postulate II, according to the way we interpreted 'continuously'.

The attitude of modern mathematicians is that one can vary the postulates of Euclid at will, constructing as many different geometries as one wishes. In the 19th century, this was a radical idea. People thought of Euclid's axioms as necessary truths about space, and hence truths which underlay the whole of astronomy and physics.

A modern physicist uses whichever geometry suits his purposes. According to the general theory of relativity, space-time is a four-dimensional Riemannian geometry, but with its curvature varying from place to place, depending on the local density of matter. The sum of the angles of a triangle might be two right angles (as in Euclidean geometry) if one was in a vacuum; however, if matter were present, the angle sum would differ from two right angles (as in non-Euclidean geometry), on account of the bending of light rays under the gravitational influence of that matter. Ideas such as these would have amazed mathematicians living in the early part of the 19th century.

We have seen in the last chapter that Euclid's postulates were not really adequate to describe the system he had in mind. Surprisingly, it was only in 1899 that Hilbert gave a completely adequate axiomatic description of three-dimensional Euclidean space. Since Hilbert required 21 postulates, or 'axioms' as he preferred to call them, we shall only state some of them here

to give the flavour of his work.

Hilbert deals with the following undefined concepts: point, line, plane, incidence (between points and lines, between points and planes, between lines and planes), order (a ternary relation of 'betweenness' for three collinear points) and congruence (a binary relation between 'segments', which are themselves defined in terms of betweenness).

He lists seven axioms of incidence. For example, the first two can be combined to say this:

'Given two distinct points A and B there is a unique line a such that A lies on a and B lies on a .'

He lists five axioms of order. For example, the first of these says this:

'If B is between A and C then B is between C and A .'

More significant is his fifth axiom of order:

'If A , B , and C are three non-collinear points, and if a is a line which meets the segment AB , then a also meets the segment AC or the segment BC .'

He lists six axioms of congruence; for example the second one says

'If $AB \equiv A'B'$ and $AB \equiv A''B''$ then $A'B' \equiv A''B''$.'

He lists two axioms of continuity. The first of these is the so-called axiom of Archimedes: 'If e and f are geometric quantities and $e \neq 0$, then there is a natural number n such that $ne > f$.'

The second is his controversial axiom of completeness. He only reluctantly added it to the French translation of his lecture notes when it became apparent that otherwise one still could not deduce Euclid's Proposition 1. It was later simplified by Bernays as follows:

'No points can be added to a straight line so that all other postulates remain valid.'

Exercises

1. How does one use the parallel postulate to show, in Euclidean geometry, that every triangle has a circumcircle?
2. Show that, in the Poincaré model, there is exactly one line through any two distinct points. That is, prove, in Euclidean geometry, that, given any two points in the interior of a circle, there is exactly one other circle which goes through those points and is orthogonal to the first circle.
3. 'The true geometry is the one which is the simplest and most beautiful.' Write a short essay on this statement, saying something about the relation between the simple, the beautiful and the true.

Alexandria from 300 BC to 200 BC

The school of mathematics established by Euclid in Alexandria produced some first rate mathematicians in the third century BC. Among them were the following:

- Aristarchus of Samos, 310 – 250 BC,
- Archimedes of Syracuse, 287 – 212 BC,
- Apollonius of Perga, 260 – 190 BC,
- Eratosthenes of Cyrene, 275 – 195 BC.

Aristarchus came from Samos, the same Greek island Pythagoras came from. He gave an interesting application of mathematics to astronomy. Let SEM be the triangle whose vertices are the sun (S), the earth (E) and the moon (M) (Figure 18.1). Aristarchus noted that when the moon is at its first quarter, the angle SME is a right angle. This is why we see exactly half of the part of the moon's surface that faces the earth. When the moon is in its first quarter, one can see the sun and the moon together in the sky, at the same time. Thus Aristarchus was able to measure the angle SEM . He found it to be $29/30$ of a right angle. (A more accurate value is 0.9981 of a right angle.) Constructing a right triangle with an acute angle of $29/30$ of a right angle — there is a ruler and compass construction for this — Aristarchus found that the ratio of its short side to its hypotenuse is about $1/19$. He concluded that the distance from the earth to the sun is about 19 times greater than the distance from the earth to the moon. Had his

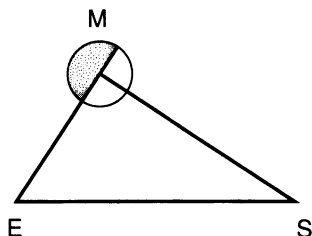


FIGURE 18.1. Relative distances of the sun and moon

measurement of the angle SEM been correct, he would have found that SE is about 400 times SM .

His calculation would have been easier had he used trigonometry, which was only developed a century later. If $\angle SEM = 87^\circ$, then $\angle ESM = 3^\circ = \pi/60$ radians, hence $EM/ES = \sin(\pi/60) \approx \pi/60 \approx 1/19$, as the sine of a small angle is approximately equal to that angle when expressed in radian measure.

Since the apparent sizes of the sun and the moon are approximately equal, as is seen during a solar eclipse, their actual diameters are in the same ratio as their distance from the earth.

By looking at the shadow cast by the earth upon the moon during a lunar eclipse one may also compare the size of the moon with that of the earth. (Since the sun is far away, the size of the earth is approximately the same as that of its shadow.) Aristarchus found

$$\frac{\text{diameter of the earth}}{\text{diameter of the moon}} \approx 7.$$

The actual figure is about 4. According to Plutarch, Aristarchus also proposed the hypothesis that ‘the earth moves in an oblique circle about the sun at the same time as it turns around its axis’. It seems that Copernicus suppressed his acquaintance with the work of Aristarchus!

Although Archimedes is assumed to have studied in Alexandria, his productive life was spent in Syracuse. We shall leave him to Chapter 19.

Apollonius (260–190 BC) came from Perga in the south of what is now Turkey. He wrote a treatise on conics which contained 400 propositions. These were arranged in eight books, four of which survived in the original Greek, and three of which survived in Arabic translation. We do not read this treatise anymore, because we feel we can do the same things more easily using analytic geometry.

According to the modern definition of a *conic section*, it is the set of all points P in the plane such that P ’s distance from a fixed point, called the *focus*, bears a constant ratio to its distance from a fixed line, called

the *directrix*. This ratio is called the *eccentricity*. Apollonius did not give this definition and it is doubtful whether he was aware of the eccentricity, which is used to classify conic sections as follows: ellipses have eccentricity < 1 (in particular, a circle has eccentricity 0), parabolas have eccentricity $= 1$, and hyperbolas have eccentricity > 1 .

Apollonius also wrote a treatise on 'Tangencies' in which he showed how to give a ruler and compass construction for a circle tangent to three given circles.

Eratosthenes of Cyrene (in North Africa) became the chief librarian at Alexandria. He was interested in many things: philosophy, poetry, history, philology, geography, astronomy and mathematics. We have already mentioned his sieve for constructing the list of primes. Eratosthenes also invented the Julian calendar, with every fourth year containing an extra day, and he calculated the size of the earth. Perhaps the reason his students called him 'beta' (the second letter of the Greek alphabet) was that, although he studied many different things, he never considered himself the leading expert in any one field. It is reported that, in his old age, Eratosthenes went blind and committed suicide by starvation.

Eratosthenes's greatest achievement was the measurement of the circumference of the earth. Eratosthenes correctly assumed that, since the sun is so far from the earth, those of its rays which hit the earth can be regarded as parallel. (Here he used the result of Aristarchus.) Eratosthenes knew that Syene (present day Aswan) is almost exactly on the Tropic of Cancer, that is, at noon on midsummer's day (June 21), the sun is directly overhead, as could be witnessed from the bottom of a well. Eratosthenes observed that at Alexandria, at noon on midsummer's day, the sun was $360^\circ/50$ from the point directly overhead. He argued that this same angle was subtended at the center of the earth by the arc joining Alexandria to Syene, which is due south of Alexandria. According to Euclid (theorem VI 33), the length of an arc of a circle is proportional to the angle it subtends at the center. So all Eratosthenes had to do was to measure the distance from Alexandria to Syene. This he found to be 5,000 stadia, a stadium being the length of the famous Olympic track. Eratosthenes concluded that the circumference of the earth is to 5,000 stadia as 360° is to $360^\circ/50$, and hence the circumference is $5,000 \times 50 = 250,000$ stadia (Figure 18.2). As the Olympic stadium is about 180 meters, this would make the circumference of the earth about 45,000 kilometers.

Eratosthenes's calculation of the circumference of the earth is remarkably accurate. The correct value is almost exactly 40,000 km; in fact, the *kilometer* was originally defined as $1/40,000$ of the circumference of the earth. Had Columbus known this, he might never have set out on his journey or called the inhabitants of the New World 'Indians'.

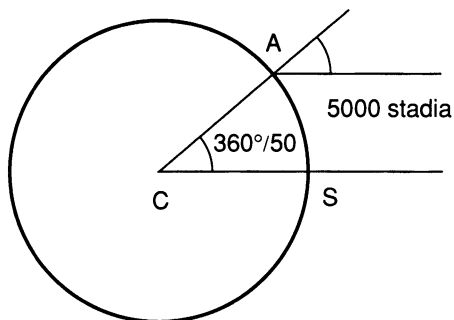


FIGURE 18.2. Circumference of the Earth in stadia

Exercises

1. Obtain the equation of a conic section with focus $(a, 0)$, directrix the y axis and eccentricity e .
2. Suppose you know the actual size of the moon. What is a simple way of finding its distance from the earth — without using anything Eratosthenes could not have used.

Archimedes

Archimedes (287–212 BC) was the greatest applied mathematician and physicist before Newton. Many stories are told about him. One story relates that, while he was taking a bath, Archimedes suddenly discovered a simple way of determining the ratio of gold to silver in a gold-silver alloy. Elated by his discovery, he leapt from the bath, and ran through the streets of Syracuse, shouting ‘eureka!’ (which means ‘I have found it!’). Unfortunately, he had forgotten to get dressed.

It was no accident that Archimedes made his discovery in a bath. Suppose you have an object made of gold and silver weighing m ounces. Suppose you wish to determine the number x of ounces of gold which the goldsmith has put into it. If g is the density of gold and s is the density of silver, the volume of the object is $(x/g) + (m - x)/s$. What Archimedes realized was that, by immersing the object in a rectangular bath tub, and observing the increase in water level, you can easily determine its volume v . Solving for x in the equation

$$v = (x/g) + (m - x)/s,$$

you obtain the mass of gold in the object. Thanks to his mathematics, Archimedes was able to tell his friend, King Hieron of Syracuse, whether the goldsmith had cheated the king by charging him for pure gold, while, in fact, using a certain percentage of silver for his crown. Newton, the first modern physicist to surpass Archimedes, centuries later, was to perform a similar task in unmasking counterfeiters.

When Syracuse was besieged by a Roman army, Archimedes constructed various machines to help defend his city. As well as catapults and cross-

bows, Archimedes designed devices which dropped huge stones on the Roman ships. He even constructed a crane which could lift a ship from the water, and drop it back in, stern-first. When Syracuse finally fell (212 BC), the Roman general Marcellus gave orders to bring Archimedes to him unharmed. These were not obeyed. Archimedes, it seems, was slain by an unknown soldier. There are various accounts of this story; the best known is that by Plutarch in his biography of Marcellus. Who today would know of Marcellus were it not for Archimedes?

Archimedes wrote on many subjects: the circle, the parabola, the spiral, the sphere, the cylinder, arithmetic, mechanics, statics and hydrostatics. One of his more interesting books, the *Method*, was rediscovered only in 1906.

By considering a regular 96-gon inscribed in a circle, Archimedes showed that $\pi < 3\frac{1}{7}$. By considering a regular 96-gon circumscribing a circle, he showed that $3\frac{10}{71} < \pi$. He was aware that one could calculate π to any desired accuracy by letting the number of sides of the regular polygon tend to infinity.

He proved that the area of a circle is πr^2 and that the volume of a sphere is $\frac{4}{3}\pi r^3$. He knew how to calculate the area bounded by a parabola and a chord, the area of a sector of a spiral, the volume of an ellipsoid of revolution, the volume of a segment of a sphere, the centroid of a hemisphere and, perhaps most remarkably, the volume common to two equal right circular cylinders intersecting at right angles. All these are calculus problems and, indeed, Archimedes was using what we would now call the technique of 'integration'.

As an example of Archimedes's mathematics, let us see how he proved that the area of a circle is πr^2 . He started with the following assumptions and theorems:

1. Circles and circle segments have areas.
2. The area of a set of pairwise disjoint triangles and circle segments equals the sum of the areas of those triangles and circle segments; thus if we dissect a circle into triangles and circle segments, the area of the circle is the sum of the areas of the triangles and circle segments into which it has been dissected; also the area of the circle is greater than the sum of the areas of any proper subset of those triangles and circle segments.
3. Given any circle, there is a straight line segment which is longer than the perimeter of any polygon inscribed in the circle, and shorter than the perimeter of any polygon circumscribing the circle; this is the 'circumference' of the circle. (The first person to give the name π to the circumference of the circle with unit diameter was not Archimedes, but William Jones, in 1706.)

4. Given any areas e and f , there is a natural number m such that $me > f$; (this assumption is found at the beginning of Archimedes's 'On the Sphere and the Cylinder', so it is often called the 'Axiom of Archimedes'; however, it is also found at the beginning of Book V of Euclid's *Elements* and, even earlier, at 266b in Aristotle's *Physics*).
5. A regular 2^n -gon inscribed in a circle takes up more than $1 - 1/2^{n-1}$ of its area; a regular 2^n -gon circumscribed about a circle has an area less than $1 + 1/2^{n-2}$ of that of the circle.
6. The area of a circle is proportional to its diameter squared (see Euclid's *Elements* XII 2).

Using these assumptions and theorems, Archimedes derived the formula for the area of a circle by first obtaining two contradictions:

(A) Suppose the circle has area greater than that of a right triangle T whose legs equal the radius and circumference of the circle.

By (4) and (5) we can find a natural number n such that

$$\text{circle area} - \text{inscribed regular } 2^n\text{-gon area} < \text{circle area} - \text{area of } T$$

and hence

$$\text{area of } T < 2^n\text{-gon area.}$$

Let AB be a side of the inscribed regular 2^n -gon, and ON a perpendicular from the center O of the circle to AB (with N being the midpoint of AB). Then ON is less than the radius of the circle. Using (3), we have

$$\begin{aligned} 2^n\text{-gon area} &= 2^n(\tfrac{1}{2}AB \cdot ON) \\ &= \tfrac{1}{2}(2^n AB)ON \\ &< \tfrac{1}{2}\text{circumference} \times \text{radius} \\ &= \text{area of } T. \end{aligned}$$

Contradiction. Thus (A) must be rejected.

(B) Suppose the circle has area less than T .

By (4) and (5) there is a natural number n such that the area of $T >$ circumscribed regular 2^n -gon area. However, if AB is a side of the circumscribing regular 2^n -gon, then, by (3),

$$\begin{aligned} 2^n\text{-gon area} &= 2^n(\tfrac{1}{2}AB \times \text{circle radius}) \\ &> \tfrac{1}{2}\text{circumference} \times \text{radius} \\ &= \text{area of } T. \end{aligned}$$

Contradiction. Thus (B) must be rejected.

Since (A) and (B) must be rejected, it follows from Aristotle's 'Law of the Excluded Middle' that (C) the area of the circle equals that of a right triangle whose legs equal the radius and circumference of that circle. In other words,

$$\text{circle area} = \frac{1}{2} \text{ circumference} \times \text{radius}.$$

The expression on the right-hand side is of course πr^2 (see (6) above).

Archimedes's proof of this formula was the culmination of about two hundred years of previous work on the circle, beginning with Antiphon (425 BC).

Archimedes, working in Syracuse, would communicate his results to the mathematicians back in Alexandria. He became annoyed when, suspiciously often, they claimed that they had made the same discoveries. To fool them, Archimedes included some false results in a book on the sphere and the cylinder, but history does not reveal the outcome. To challenge the mathematicians at Alexandria, Archimedes posed the following problem (which we have 'translated' into modern algebraic notation).

The sungod had a herd of cattle consisting of w white bulls, g grey bulls, b brown bulls and s spotted bulls, as well as w' , g' , b' and s' cows of matching colours. What was the total number of bulls and cows if

$$s = w - 5g/6 = g - 9b/20 = b - 13w/42,$$

$$w' = 7(g + g')/12, \quad g' = 9(b + b')/20,$$

$$b' = 11(s + s')/30, \quad s' = 13(w + w')/42,$$

and $w + g$ is a square and $b + s$ is a triangular number?

It is a curious triumph of tradition that Archimedes used the Egyptian method for representing the fractions appearing in this problem as sums of reciprocals of positive integers. Aside from the last two restrictions, we have here seven equations in eight unknowns, which cannot be solved by algebraic methods alone. However, if one is looking for positive integer solutions, such problems are called 'Diophantine', after the mathematician Diophantus, who will appear about 500 years after Archimedes. We sketch how a solution may proceed, though the reader will have to fill in many details.

From the equations

$$s = w - 5g/6 = g - 9b/20 = b - 13w/42$$

we find that, for some positive integer m ,

$$w = 2226m, \quad g = 1602m, \quad b = 1580m, \quad s = 891m.$$

From the next four equations, we find that there is some natural number k such that

$$m = 4657k, \quad w' = 7,206,360k, \quad g' = 4,893,246k,$$

$$b' = 3,515,820k, \quad s' = 5,439,213k.$$

If $w + g$ is a square then $4(957)(4657)k$ is a square. Since 4657 is prime and 957 is the product of two distinct primes, it is easy to show that this occurs only if k has the form $(957)(4657)t^2$, where t is an integer. If $b + s$ is triangular, then $2471m$ has the form $\frac{1}{2}n(n+1)$. Thus $8(2471)(4657)(957)(4657)t^2$ has the form $4n^2 + 4n$. In other words, to find an integer solution to Archimedes's Cattle Problem, we have to find an integer solution to

$$(2n+1)^2 - 8(2471)(957)(4657^2)t^2 = 1.$$

The mathematicians at Alexandria were not able to solve this problem. Indeed, it was not solved until 1965, when H. C. Williams, R. A. German and C. R. Zarnke used a computer to generate the 206,545 digit answer. The answer was published for the first time in 1980–1981. The reader can find the full 206,545 digit answer printed in Harry L. Nelson's article 'A Solution to Archimedes' Cattle Problem', *Journal of Recreational Mathematics* 13, pp. 164–176.

It is impossible to do full justice here to Archimedes's important contributions to physics. Let us only mention that he developed the theory of the lever and investigated the properties of floating bodies.

Exercises

1. Let a and b be positive real numbers. Archimedes proved that $x^3 - ax^2 + (4/9)a^2b$ has a positive root if and only if $a > 3b$. Do the same.
2. If, in a cube of side 1, two cylinders, each of diameter 1, are constructed so that their axes are perpendicular, show that the volume common to these cylinders is $2/3$.
3. Prove that a regular 2^n -gon circumscribed about a circle has an area less than $1 + 1/2^{n-2}$ of that of the circle.
4. Find the least positive integer solution of $x^2 - 5y^2 = 1$.
5. Let B be a point in straight line segment AC . Construct three semi-circles with diameters AB , BC and AC , all on the same side of AC . The area which is in the semi-circle on AC but not in either of the two smaller semi-circles is called the 'arbelos' (or 'shoe-maker's knife'). Archimedes found the area of the arbelos, in terms of AB and BC . Do the same.

6. Justify the existence of the numbers m and k in the above solution of the Cattle Problem.
7. If Archimedes were alive today, would he have a moral obligation to help his country design nuclear weapons or would he have a moral obligation not to help them design nuclear weapons? Support your answer with reasons related to the role of science in human history.

Alexandria from 200 BC to 500 AD

In this chapter we discuss the more important mathematicians who worked in Alexandria after 200 BC:

- Hipparchus of Nicea, born about 180 BC,
- Heron of Alexandria, about 60 AD,
- Menelaus of Alexandria, about 100 AD,
- Ptolemy of Alexandria, died in 168 AD,
- Diophantus, about 250 AD,
- Pappus, about 320 AD.

Less significant as mathematicians, but nonetheless important in the history of the subject are

- Nicomachus of Gerasa, about 100 AD,
- Hypatia, died in 415 AD,
- Proclus, 410 – 485 AD,
- Boethius, 475 – 524 AD.

Hipparchus came from Nicea, a town near present day Istanbul, which was to be the site of the great Church Council of 325 AD. Hipparchus made many contributions to astronomy. He calculated the duration of the year

to within 6 minutes, the angle between the ecliptic and the equator, the annual precession of the equinoxes, the lunar parallax, the eccentricity of the solar orbit, etc. He knew that the moon moves only approximately in a circle with center at the earth, a better approximation being an 'epicycle'. The same was true of the sun. Hipparchus suggested that epicycles of higher orders were necessary to describe the motions of the planets.

In mathematics his great contribution was the founding of trigonometry. He drew up a table giving for each angle with vertex at the center of a circle of radius 1 the length of the chord it cuts off in the circle. For example, suppose $\angle AOB$ is 30° , with O the center of the circle, and $OA = OB = 1$ two radii of the circle. Then the chord in question is the segment AB . This has length $31.06/60$, so, in the table of Hipparchus, we would find

$$\text{chord}(30^\circ) = 31.06/60.$$

In modern terms, $\text{chord}(x) = 2 \sin(x/2)$.

To construct this table, Hipparchus made use of formulas which we would express as follows:

$$\sin(x \pm y) = \sin x \cos y \pm \cos x \sin y$$

and

$$2 \sin^2(x/2) = 1 - \cos x.$$

Heron is known for the formula, probably discovered by Archimedes, which expresses the area of a triangle in terms of the lengths a, b and c of its sides. If $s = (a + b + c)/2$, this formula gives the area of the triangle as

$$\sqrt{s(s-a)(s-b)(s-c)}.$$

We shall meet this formula again when we discuss the mathematics of India.

Menelaus was the first to study spherical trigonometry. He is also known for the following theorem, which is found in his *Spherica*.

Menelaus's Theorem:

Let ABC be a triangle. Suppose D is on the line through B and C , E is on the line through A and C , and F is on the line through A and B . Suppose that either two or none of D, E, F are on sides of the triangle. Then D, E, F are collinear if and only if $BD \cdot CE \cdot AF = CD \cdot AE \cdot BF$.

Proof: Suppose D, E, F are collinear, in line l . We may suppose l does not pass through A, B or C . Let A', B', C' be points in l such that AA', BB', CC' are all perpendicular to l . Then $BD/CD = BB'/CC', CE/AE = CC'/AA'$

1	α	10	ι	100	ρ
2	β	20	κ	200	σ
3	γ	30	λ	300	τ
4	δ	40	μ	400	υ
5	ϵ	50	ν	500	ϕ
6	f	60	ξ	600	χ
7	ζ	70	o	700	ψ
8	η	80	π	800	ω
9	θ	90	q	900	

TABLE 20.1. Ptolomaic notation

and $AF/BF = AA'/BB'$. Multiplying the three equations together we obtain the result.

The converse is easy.

As his *Tetrabiblos* shows, Ptolemy was a keen believer in the superstition called astrology. In spite of this, he did for astronomy what Euclid had done for geometry and arithmetic: he wrote the definitive textbook, known by its Arabic name, the ‘Almagest’ or ‘Greatest’. Like Hipparchus, Ptolemy gave a table of chords.

Book I of the ‘Almagest’ contains ‘Ptolemy’s Theorem’: in a cyclic quadrilateral, the product of the diagonals is equal to the sum of the products of the two pairs of opposite sides. He used Greek letters to denote numbers. Curiously, however, he retained two Phoenician letters, corresponding to Latin f and q , which had actually disappeared in Greek. He also added one symbol at the end, giving him a total of 27 symbols, which allowed him to represent the numbers 1 to 9, 10 to 90 and 100 to 900. He also made use of a small circle to denote zero. In his tables, he employed the Babylonian system to denote not only angles, as we still do, but also lengths as had Hipparchus before him. Thus he wrote

$$120^{\circ}0'0'' = \rho\kappa|\circ|\circ,$$

$$\text{chord } 1^{\circ} = 1^{\circ}2'50'' = \alpha|\beta|\nu.$$

So he took the radius to be 1.

His underlying decimal notation was based on the alphabetic code in Table 20.1. We have substituted the Latin letters for 6 and 90 and omitted the symbol for 900.

In 250 AD, in Rome, Plotinus was teaching his version of Platonism. At the same time, in Alexandria, Diophantus was writing the *Arithmetica*. This originally contained 13 books. Until 1973 we had only six of these,

but then three more were discovered in an Arabic translation going back to the ninth century. (See Jacques Sesiano, *Books IV to VII of Diophantus's Arithmetica*.)

These 13 books consisted of solutions to algebraic problems. The solutions are all rational numbers. Some of the equations or systems of equations are indeterminate and often there is more than one rational solution. Diophantus, however, is usually content to give just one solution. Note that what we now call a 'Diophantine equation' is one whose unknowns are not just rational but integers. Diophantus, however, accepted any rational solution.

As an example, let us consider problem 9 of Book II. The problem is 'to divide a given number which is the sum of two squares into two other squares.' That is, given rationals a and b , find a nontrivial rational solution of

$$x^2 + y^2 = a^2 + b^2.$$

Diophantus takes the special case where $a = 2$ and $b = 3$, but his solution is easily generalized. He writes:

Take $(x + 2)^2$ as the first square and $(mx - 3)^2$ as the second (where m is an integer), say $(2x - 3)^2$. Therefore $(x^2 + 4x + 4) + (4x^2 + 9 - 12x) = 13$, or $5x^2 + 13 - 8x = 13$. Therefore $x = 8/5$, and the required squares are $324/25$ and $1/25$.

Note that, to get the general solution, m should be any rational number.

An interesting question is whether Diophantus was aware of the algebraic rules that lay behind many of his solutions. In Book III, problem 19, Diophantus notes that 65 is a sum of two squares in two ways since 65 'is the product of 13 and 5, each of which numbers is the sum of two squares'. From this we can deduce that he knew that the product of two integers, each of which is a sum of two squares, is itself a sum of two squares, and in two ways. Did Diophantus also know the stronger proposition that

$$(a^2 + b^2)(c^2 + d^2) = (ac \mp bd)^2 + (ad \pm bc)^2 ?$$

Basing himself just on the remark which we have quoted from Book III, problem 19, T. L. Heath conjectured that Diophantus did know this identity (see page 105 in Heath's translation of the *Arithmetica*). The person who first published the algebraic identity, however, was Abu Jafar al-Khazin (950 AD), and, later, Fibonacci gave it in his *Liber Quadratorum* (1225 AD).

Diophantus was the first to make systematic use of a symbolic notation for algebraic expressions. He denoted $+$ by juxtaposition, $-$ by the symbol λ and $=$ by ι . He wrote

$$\begin{aligned} K^v &\text{ for } x^3, \\ \Delta^v &\text{ for } x^2, \end{aligned}$$

$$\varsigma \text{ for } x^1 \text{ } (\varsigma\varsigma \text{ for the plural}),$$

$$\overset{\circ}{M} \text{ for } x^0.$$

For example,

$$\Delta^v \overline{\gamma} \overset{\circ}{M} \overline{\iota\beta}$$

stands for $3x^2 + 12$, while

$$K^v \overline{\alpha\varsigma\varsigma\eta} \wedge \Delta^v \overline{\epsilon} \overset{\circ}{M} \overline{\alpha} \wr \varsigma \overline{\alpha}$$

represents the equation

$$(x^3 + 8x) - (5x^2 + 1) = x.$$

In 320 AD, the Roman Empire had its first Christian emperor, Constantine, after whom the Eastern Roman capital was named Constantinople. In Alexandria, Athanasius was defending the divinity of Jesus against Arius, who asserted that Jesus was *like* God, but not *equal* to him. (The difference between these two words in Greek consisted of one letter, the Greek letter iota. We have preserved this difference in Mathematics, when we distinguish between ‘homeomorphism’ and ‘homomorphism’.)

Meanwhile in Alexandria, Pappus was writing his encyclopaedic *Collection* of earlier mathematical works. The school of mathematics had declined, and Pappus was its last, lone member.

The ‘Theorem of Pappus’ appears in Book VII of the *Collection*. It is far more important than Pappus realized. It expresses the commutativity of multiplication and is fundamental to projective geometry. Hilbert made use of it as a key theorem in his presentation of Euclidean geometry. The theorem of Pappus can be proved with the help of the theorem of Menelaus as follows.

Theorem of Pappus: Given points ABC on one line, $A'B'C'$ on another, the three points of intersection $P = BC' \cap CB'$, $Q = AB' \cap BA'$ and $R = CA' \cap AC'$ are collinear (Figure 20.1).

In stating this result, we have assumed that BC' and CB' etc., are not parallel. We shall also assume that ABC and $A'B'C'$ meet at a point X and that none of the other lines in the diagram are parallel.

Proof: Let $A'B \cap B'C = U$, $AC' \cap A'B = V$ and $B'C \cap AC' = W$. We apply Menelaus’s Theorem five times to the triangle UVW . Since $A'CR$, $BC'P$, $AB'Q$, $A'B'C'$, and ABC are all collinear, we have

$$VR \cdot WC \cdot UA' = RW \cdot CU \cdot A'V,$$

$$VC' \cdot WP \cdot UB = C'W \cdot PU \cdot BV,$$

$$VA \cdot WB' \cdot UQ = AW \cdot B'U \cdot QV,$$

$$VC' \cdot WB' \cdot UA' = C'W \cdot B'U \cdot A'V,$$

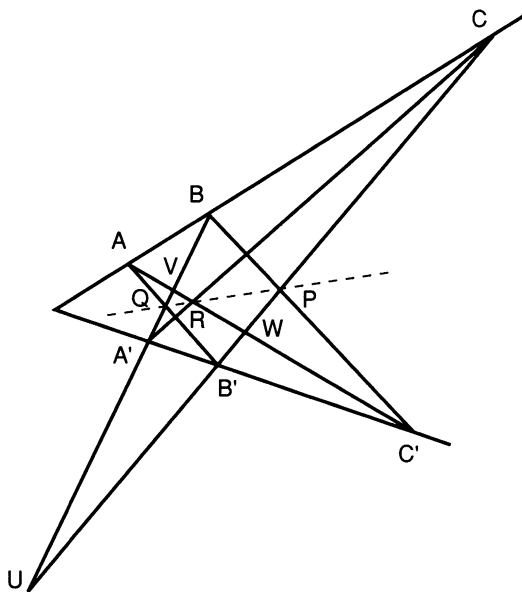


FIGURE 20.1. Pappus's Theorem

$$VA \cdot WC \cdot UB = AW \cdot CU \cdot BV.$$

Multiplying the first three equations and dividing by the product of the last two, we obtain

$$VR \cdot WP \cdot UQ = RW \cdot PU \cdot QV.$$

The result now follows by another application of Menelaus's Theorem. Note that the argument makes use of the commutativity of multiplication.

When CD and EF are parallel, the proof proceeds in a similar fashion. However, one uses, not Menelaus's Theorem, but the theory of similar triangles. We leave the details to the reader.

Among the minor mathematicians of this era was Nicomachus (100 AD) from Palestine. He was a Pythagorean, and he published a book on number theory, which is the basis for many of our speculations about the nature of 'Pythagorean' mathematics.

Hypatia (d. 415) was the daughter of Theon of Alexandria, who had put out an edition of Euclid's *Elements* and a commentary on Ptolemy's *Almagest*. Hypatia wrote commentaries on Apollonius and Diophantus.

According to Socrates Scholasticus (380–450 AD) in Chapter 15 of Book VII of his *History of the Church*, Hypatia was murdered by a mob in the course of an anti-pagan riot. This tragedy is sometimes blamed on the Christian bishop, Cyril, but there is no evidence to support this accusation. The 19th century author C. Kingsley wrote a fascinating historical novel *Hypatia*, which makes this story come to life.

Another minor mathematician at this time was Proclus. He studied in Alexandria and then worked in Athens. He wrote a commentary on the first book of Euclid, which contains valuable information about the history of Greek mathematics.

Boethius studied in Athens but lived in Rome. He is most famous for his *De Consolatione Philosophicae*, which he wrote in prison, about 525 AD. His *Arithmetic* and *Geometry* were standard textbooks in the Middle Ages. Unfortunately, they contained much less mathematics than the *Elements*.

In 529 AD, the emperor Justinian closed the pagan schools of philosophy at Athens. The 'Dark Ages' of Europe had begun.

It is interesting that we do not know of many mathematicians of the period (50–500 AD) converting to Christianity. It seems that the Academy at Athens and the University at Alexandria both rejected the new religion. An interesting dialogue between reason and faith might have taken place, but, as it turned out, it was only in the later Middle Ages that thinkers, such as Aquinas (1250 AD), advanced philosophies that were influenced by the *Elements* as well as by the Bible.

When the Arabs conquered Alexandria in 641 AD, there probably was not much left of the famous library. However, according to an often repeated story, dating back to Moslem sources in the 13th century, it was the Arabs who destroyed it. Their commander, Amru, was willing to spare the library, but was dissuaded by Caliph Omar I, who argued thus: 'If the books of the Greeks confirm what is written in the Koran, they are superfluous; if they contradict the Koran, they are dangerous. In either case they should be destroyed.' The story continues, saying that the books served to heat the furnaces of the public bathhouses for six months! An almost identical story is told about another library in Persia. What was the origin of these stories? According to Bernard Lewis, in a Letter to the Editor, *The New York Review of Books* 37, Number 14 (October 27, 1990), they were invented to justify the destruction of a completely different library, a collection of Fatimid books, deemed to be heretical by the orthodox Sunni Sultan Saladin in the 12th century.

Exercises

1. Prove Heron's formula.
2. Give the details of the converse in the proof of the Theorem of Menelaus.
3. Find all the integer solutions of $x^2 + y^2 = z^2$ by using the same technique as Diophantus did in problem 9 of Book II of the *Arithmetica*. (Hint: the given equation has a solution in the integers if and only if $x'^2 + y'^2 = 1^2 + 0^2$ has a corresponding solution in rationals.)

4. Find an infinite family of rational solutions to $x^2 + y^2 = z^2$ (*Arithmetica* IV 1 (Sesiano)).
5. Prove a version of the Theorem of Pappus in case BC' and AB' are parallel.

Mathematics in China and India

Not much is known about the development of mathematics in China before contact with the West was established. The ‘Arithmetic in Nine Sections’ (‘Chiu Chang Suan Shu’) was written before 200 AD. Like the Rhind Papyrus, it is a list of problems and solutions. Chapter 8 shows how to solve n linear equations in n unknowns, using a method which is essentially the same as Gaussian elimination. One system which is solved is the following:

$$\begin{aligned} 3x + 2y + z &= 39, \\ 2x + 3y + z &= 34, \\ x + 2y + 3z &= 26. \end{aligned}$$

The Chinese interest in systems of linear equations was perhaps linked to their interest in magic squares. The square

$$\begin{array}{ccc} 4 & 9 & 2 \\ 3 & 5 & 7 \\ 8 & 1 & 6 \end{array}$$

was supposedly brought to humankind on the back of a tortoise from the River Lo in the days of Emperor Yu. Its ‘magic’ property is that all rows and columns and the two diagonals have the same sum.

A ‘Chinese Remainder Problem’ was solved by Sun Tsu (400 AD):

divide by 3, the remainder is 2;
 divide by 5, the remainder is 3;
 divide by 7, the remainder is 2;
 what will be the number?

Only one of the infinitely many solutions is given, but Sun Tsu's method allowed him to give as many others as he wanted.

Tsu Ch'ung Chih (475 AD) gave upper and lower bounds for π . He used the method of Archimedes, but he obtained sharper estimates. In the 'Nine Sections of Mathematics' (1247), Ch'in Chiu Shao (Qin Jiushao) found a root of

$$x^4 - 763200x^2 + 40642560000,$$

using what is in effect Horner's method (rediscovered by William Horner in 1819). In the 'Precious Mirror of the Four Elements' (1303), Chu Shih Chieh (Zhu Shijie) gives 'Pascal's Triangle' (also known in India and later rediscovered by Pascal in 1653).

In the 16th and 17th centuries, Christian missionaries from Europe entered China, introducing Western mathematics (e.g., the theory of logarithms). Today, the Chinese contribute to all branches of mathematics. They have retained, however, a strong interest in the theory of numbers. For example, in 1985, De Gang Ma gave the first elementary solution of the Diophantine equation

$$x(x+1)(2x+1) = 6y^2.$$

(See Anglin [1995], Section 4.4.)

We have already heard about Chen Jing-run and the progress he made on the Goldbach conjecture.

The earliest Indian mathematical texts we have are the *Sulvasutras* or 'rules of the cord'. They were written sometime between 500 BC and 200 AD. They contain some elementary geometry related to the construction of altars. It is observed that $12^2 + 35^2 = 37^2$ and, presumably, the general method for constructing such 'Pythagorean triples' was known.

Mathematics in India turns up in unexpected places. Around 800 BC a certain Pingala wrote a book on the *Science of verse meters*. Syllables in Vedic poetry are distinguished to be either 'light' or 'heavy'; the former are assigned the binary digit 1, the latter the binary digit 0. Thus, with each line of verse is associated a number written in binary scale, which is then converted into decimal scale. Conversion in the converse direction is also discussed. (See van Nooten [1993].)

Four noteworthy Indian mathematicians were:

- Aryabhata the Elder, born in 476 AD,
- Brahmagupta, flourished in 628 AD,
- Mahavira, lived about 850 AD,
- Bhaskara, 1114–1185 AD.

In their work, ancient Indian traditions and Alexandrian mathematics seem to flow together. Yet, unlike the ancient Greeks, these mathematicians

did not include many proofs in their books. Also, unlike the ancient Greeks, they wrote their mathematics books in verse! For example, Bhaskara gives the following problem:

The square root of half the number of bees in a swarm
Has flown out upon a jasmine bush;
Eight ninths of the swarm has remained behind;
And a female bee flies about a male who is buzzing inside a
lotus flower;
In the night, allured by the flower's sweet odour, he went inside
it
And now he is trapped!
Tell me, most enchanting lady, the number of bees.

This is certainly a poetic way of asking for the solution of

$$\sqrt{x/2} + 8x/9 + 2 = x.$$

Aryabhata wrote on arithmetic and algebra, on plane and spherical trigonometry, and on astronomy. He knew the formulas for the sum $1^k + 2^k + \dots + n^k$ for $k = 1, 2$ and 3 . He solved quadratic and indeterminate equations. He was one of the first to make use of the sine function, defining $\sin A$ as $(\text{chord } 2A)/2$, thus helping to simplify the addition formulas and giving trigonometry its modern form. He also produced a quite accurate sine table.

Brahmagupta (b. 598) studied the Diophantine equation $x^2 - Ry^2 = 1$ (where R is a given, positive nonsquare integer). This equation is mistakenly called 'Pell's equation', after John Pell (1610–1685), who actually had nothing to do with it. A complete solution to this equation was given by Lagrange in 1767.

Brahmagupta calculated the surface and volume of pyramids and cones, using $\sqrt{10}$ as an approximation for π . Most importantly, he was the first person to give a systematic presentation of the rules for working with zero and negative numbers, and he may have been the first to introduce negative numbers into the number system. According to *The Treasury of Mathematics*, edited by Henrietta O. Midonick, Brahmagupta's *Brahmasphuṭa Siddhanta* tells us that positive divided by positive, or negative divided by negative, is positive, whereas positive divided by negative, or negative divided by positive is negative. It seems, however, that Brahmagupta had some trouble dividing by zero.

One of Brahmagupta's more technical achievements was his discovery of the formula for the area of a cyclic quadrilateral with sides a, b, c and d :

$$\sqrt{(s-a)(s-b)(s-c)(s-d)},$$

where $s = (a + b + c + d)/2$, thus generalizing Heron's formula. He did not give a proof for this formula, and it seems that a careless copyist omitted the word 'cyclic'.

Following Brahmagupta, Mahavira discusses the properties of zero, stating that

$$a + 0 = a, a - 0 = a, a \cdot 0 = 0.$$

He does not quite know what to do with $a/0$, apparently believing that it equals a , unless we again blame a copyist. He also talks about negative numbers and asserts that they are not squares. He tries to give some social significance to his equations. When posing the equation

$$x/4 + 2\sqrt{x} + 15 = x,$$

he writes:

One fourth of a herd of camels was seen in a forest. Twice the square root of that herd had moved on to the mountain slope. Three times five camels, however, were found to remain on the river bank. What is the number of that herd of camels?

Bhaskara (1114–ca. 1185) named one of his mathematics books after his daughter Lilavati. This book deals with weights and measures, the decimal notation, the operations of arithmetic, reduction of fractions to common denominators, linear and quadratic equations, arithmetic and geometric progressions, interest and discount, triangles and quadrilaterals, approximations to π , trigonometric formulas, volumes of solids, indeterminate linear equations, and combinations. It contains the earliest extant exposition of the decimal notation using the sign for zero. Concerning calculations with zero, Bhaskara gives the following rules:

$$a \pm 0 = a, 0^2 = 0, \sqrt{0} = 0, a/0 = \infty.$$

Some of his problems throw light on economic and social history. A female slave of age 16 is said to be worth eight oxen who have worked for two years; her price declines as she grows older.

In the first volume of his *History of Mathematics*, D.E. Smith relates that the *Lilavati* was translated into Persian in 1587 by someone called Fyzi. Fyzi claims that Bhaskara dedicated the book to his daughter in order to console her for remaining single. According to astrologers, there was but one lucky moment at which Lilavati might marry. Unfortunately, one of Lilavati's pearls fell into the water clock, and it stopped without anyone noticing. The lucky moment passed, and Lilavati missed her chance of an astrologically sound marriage. Bhaskara, accepting the advice of the astrologers, decided to give Lilavati a book of mathematics instead of a husband.

Perhaps the greatest Indian mathematician in the 20th century was Srinivasa Ramanujan (1887–1920), who was essentially self-taught and discovered many beautiful and ingenious formulas.

Exercises

1. Solve the system of three equations in three unknowns given in the Chiu Chang Suan Shu.
2. Find a 4 by 4 magic square, using the numbers from 1 to 16.
3. Find all the solutions to Sun Tsu's Chinese Remainder Problem.
4. Show that Ch'in Chiu Shao's polynomial had four linear factors.
5. Find a solution of De Gang Ma's equation (first posed by E. Lucas in 1875) with $x > 1$.
6. How many camels were there in Brahmagupta's problem?
7. How many bees were there in Bhaskara's problem?
8. If you know the sides of a cyclic quadrilateral, can you determine the radius of the circumscribing circle? Why? How?

Mathematics in Islamic Countries

When we speak of Arabic mathematicians, we must remember that Arabic was the common language of intellectuals in the Islamic world, just as Latin was in medieval Europe. In fact, the mathematicians may have been Turks or Persians.

During the period we are concerned with here, the intellectual center of the Arab world was Baghdad, which was founded in the 8th century and was the seat of the eastern Caliphs. It was visited by Greek and Jewish physicians from the West and by Indian and Persian scholars from the East.

Caliph al-Mamun established a 'House of Wisdom', or university, at Baghdad at the beginning of the 9th century. He ordered the translations of many Greek manuscripts, and it was thanks to his policy that many of these works were preserved. Three mathematicians associated with the House of Wisdom are

- Al-Khwarizmi, about 830,
- Thabit Ibn-Qurra (836–901),
- Al-Khayyami (Omar Khayyam) (ca. 1050–1123).

Muhammed ibn-Musa al-Khwarizmi probably came from Khorezm, which corresponds to the modern Khiva and the district surrounding it, south of the Aral Sea in central Asia. But some people deny this and say that he was merely a member of the Khwarizmi tribe, which in fact came to rule the Turko-Persian empire in 1194 AD. Our word 'algorithm' comes from a

book he wrote on the use of Indian numerals with the title: ‘Spoken has al-Khwarizmi ...’, or, in the English translation of the Latin translation: ‘Spoken has Algoritmi ...’.

Al-Khwarizmi’s most important book was the *Hisab al-jabr w’al-muqabalah*, from which we get the word ‘algebra’. The word ‘al-jabra’ means something like ‘combining’, as in combining the terms to solve an equation. The same root shows up in old Spanish as ‘algebrista’, a bone-setter, that is, one who joins together the parts of a broken bone.

Al-Khwarizmi’s ‘Algebra’ was based on the work of Brahmagupta. Although it became extremely influential, it is not as interesting from a mathematical point of view as its fame would suggest. It contains nothing that was not known to the ancient Babylonians or ancient Greeks. There are few proofs, and one of them is woefully inadequate. This is al-Khwarizmi’s ‘proof’ of the Theorem of Pythagoras, which only works if the right triangle is isosceles!

Al-Khwarizimi gives three approximations for π . None of them is supported by any reasoning, and al-Khwarizmi does not seem to think it matters which one is used. However, the book is extremely interesting as a source of sociological information. Many of the text problems deal with questions of inheritance according to Muslim religious law. (These problems appear to have been omitted in the Latin translation by Robert of Chester, which popularized al-Khwarizmi’s work in Europe.) Here is an example:

A man, in his illness before his death, makes someone a present of a slave girl, besides whom he has no property. Then he dies. The slave girl is worth 300 dirhams, and her dowry is 100 dirhams. The man to whom she has been presented cohabits with her. What is the legacy?

Here is how one is supposed to solve the problem:

$$300 - x - \frac{100}{300}x = 2x,$$

yielding $x = 90$ dirhams (Rosen [1831]).

Thabit Ibn-Qurra was an extraordinary polymath. He lived in Baghdad and was an active member of a neo-Pythagorean group called the Sabians. He wrote on politics, grammar, symbolism in Plato’s *Republic*, smallpox, the anatomy of birds, the beam balance, the salinity of seawater, the sundial, Euclid’s Parallel Postulate, cubic equations, the new crescent moon, etc.

Thabit believed there is an actual infinity (as opposed to Aristotle’s potential infinity). He did work in spherical trigonometry and in what we now would call calculus. In his *Book on the Determination of Amicable Numbers*, he gave a wholly original rule:

Let n be a positive integer greater than 1. Let $p = 3 \cdot 2^n - 1$, $q = 3 \cdot 2^{n-1} - 1$ and $r = 9 \cdot 2^{2n-1} - 1$. If p , q and r are primes, then $2^n pq$ and $2^n r$ are 'amicable' (that is, the sum of the proper divisors of each equals the other).

When $n = 2$, we get the amicable pair 220 and 284. When $n = 3$ (or any multiple of 3), r is divisible by 7, and so is not prime. However, when $n = 4$, we have another amicable pair, namely, 17296 and 18416.

Al-Khayyami wrote in Arabic on astronomy and mathematics. He revised the Julian calendar, approaching our Gregorian calendar in accuracy. In a book with the same title as that by al-Khwarizmi, he developed a geometric method for finding the positive real roots of the cubic and quartic equations. The idea was this: to solve the cubic equation

$$x^3 + ax^2 + b^2x + b^2c = 0,$$

one intersects the hyperbola $y = bc/x + b$ with the circle $(x + \frac{1}{2}(a+c))^2 + y^2 = \frac{1}{4}(a-c)^2$ and discards the point $(-c, 0)$.

As Omar Khayyam, he wrote a famous poem in Persian, the *Rubaiyat*. Its 19th century translation by Edward Fitzgerald is still a bestseller. It expounds a rather Epicurean philosophy, claiming that the most important thing is wine and the only sure thing is death:

Oh, threats of Hell and Hopes of Paradise!
 One thing at least is – *This* Life flies;
 One thing is certain and the rest is Lies;
 The Flower that once has blown for ever dies.

(See LXIII in Appendix 1 of E. FitzGerald's translation of the *Rubaiyat* (London: Bernard Quaritch, 1859).) Omar was not popular with the religious establishment of his day, and even now many Persians prefer the more mystical poetry of his fellow countrymen Hafiz.

One of Omar's contemporaries, Nizam-ul-Mulk, wrote in his autobiography that as a youth he made a mutual assistance pact with two fellow students in Naishapur, namely, Omar Khayyam and Hasan Ben Sabbah. He himself later became vizier to two Sultans of the Seljuk dynasty and was able to carry out his pledge by helping Omar to a yearly pension. Unfortunately, so he recounts, Hasan was not satisfied with the government post offered to him and ultimately became the head of a religious sect, the Ismailis. He surrounded himself by a group of fanatics, called the 'assassins', this word being derived from 'hashish'. The sect exists today, though without the practice of assassination; its leader is the Agha Khan.

Exercises

1. Show that 1184 and 1210 is an amicable pair not generated by Thabit's rule.
2. Find another amicable pair which is generated by Thabit's rule.
3. In Thabit's rule, prove that if n is a multiple of 3, then r is a multiple of 7.
4. Prove the following generalization of Thabit's rule: assuming that $p = (2^k + 1)2^{t+k} - 1$, $q = (2^k + 1)2^t - 1$ and $r = (2^k + 1)2^{2t+k} - 1$ are odd primes, then $a = 2^{t+k}pq$ and $b = 2^{t+k}r$ are amicable.
5. Let ABC be a triangle with $\angle A$ obtuse. Let B' and C' be in BC such that $\angle AB'B = \angle AC'C = \angle BAC$. Derive the following theorem of Thabit: $AB^2 + AC^2 = BC(BB' + CC')$.
6. Show that al-Khayyami's method will produce a real solution of the cubic equation in the text.

New Beginnings in Europe

Europeans only began to rouse themselves from the intellectual slumber of the Dark Ages as they came in contact with Arab civilization, mostly in Spain.

Gerbert (940–1003) had studied in Spain, where he learned the Indian numerals (but not zero). He wrote on arithmetic and geometry, which so overawed his contemporaries that they believed he had a pact with the devil. In spite of this, he became Pope and was known as Sylvester II from 999 to 1003.

One of the most difficult problems in his *Geometry* was the following: find x and y such that $x^2 + y^2 = a^2$ and $\frac{1}{2}xy = b$. This would have been an easy exercise for a Babylonian scribe!

Contemporary with Gerbert was another mathematician and churchperson, Hrotsvitha of Saxony (932–1002), who had an interest in perfect numbers.

The Englishman Adelhard of Bath (1075–1160) attended lectures at Cordova in Spain, about 1120, disguising himself as a Moslem. There he obtained a copy of Euclid's *Elements* in Arabic, which he translated into Latin. All European editions of the *Elements* were based on this translation until 1533, when the Greek original finally became available.

Abraham Ben Ezra (ca. 1095–1167), while based in Spain, travelled widely between Egypt and England. His book *Sefer ha-Mispar* explained the Hindu arithmetic, using Hebrew letters for numerals, with a zero added.

He wrote poetry of a pessimistic nature, asserting that, if he were to trade in candles, it would always be noon.

Jordanus Nemorarius (early 13th century) wrote about triangles, circles, regular polygons, Arabic numerals, primes, perfect numbers, polygonal numbers, ratios, powers and progressions. Like Diophantus, he used letters to denote the unknowns in equations.

In his *Tractatus de numeris datis*, Jordanus discusses problems of the following sort: find x and y such that $x + y = 10$ and $x^2 + y^2 = 58$. This is the sort of problem the ancient Mesopotamians were good at.

Nicole Oresme (1320–1382) was a French bishop who wrote extensively on mathematics and gave the first proof of the divergence of the harmonic series.

In the thousand years from 400 to 1400 AD there was exactly one outstanding European mathematician, namely, Leonardo of Pisa (1180–1250), who was also known as Fibonacci. A contemporary of St. Francis of Assisi, he learned his mathematics in Algeria, where his father was a custom house official. In 1202 he published the *Liber abaci*, in which he explained the Indian system of numerals and introduced the famous ‘Fibonacci sequence’

$$1, 1, 2, 3, 5, 8, 13, 21, 34, 55, \dots,$$

where each number is the sum of the preceding two. This sequence has important applications in science and in advanced number theory.

Another book by Fibonacci is his *Liber quadratorum*. This original work in indeterminate analysis gave the first proof of the identity

$$(a^2 + b^2)(c^2 + d^2) = (ac - bd)^2 + (bc + ad)^2.$$

This is equivalent to the theorem that the product of the norms of two complex numbers equals the norm of their product. (The *norm* is the square of the absolute value.) The identity implies that, if each of two integers is a sum of two squares, then so is their product.

We recall that a numerical instance of Fibonacci’s identity had already been mentioned by Diophantus (Problem 19, Book III, *Arithmetica*). An explicit statement of the identity first occurred in a commentary on Diophantus by al-Khazin. (See Anbonba [1979].)

In 1225, emperor Frederick II delayed his departure on a crusade to organize a mathematical contest. Leonardo answered all the challenges with flying colours. Two of the problems were the following:

1. Find a fraction a/b such that $(a/b)^2 \pm 5$ are both squares of fractions.
2. Solve $x^3 + 2x^2 + 10x = 20$.

Fibonacci understood negative numbers, interpreting them, in one place, as losses.

The Fibonacci numbers (1, 1, 2, 3, 5, 8, ...) proved to be so interesting that now there is a whole journal, called the *Fibonacci Quarterly*, devoted to them. Two examples of their relevance are the following.

The seeds of a sunflower head are arranged in such a way that they form two sets of arcs, emanating from the center. The number of clockwise arcs might be 21, and the number of counterclockwise arcs might be 13. In every case, these numbers are consecutive Fibonacci numbers.

One of the great feats of 20th century number theory was Matiyasevič's proof that there is no general procedure for solving Diophantine equations. We shall examine this proof in the Section on Hilbert's Tenth Problem in Part II, Chapter 21. A crucial element in this proof was Matiyasvič's use of the Fibonacci numbers.

To conclude this section, we derive the formula for the n th Fibonacci number u_n , where $u_0 = 0$ (which we include for convenience), $u_1 = 1$ and $u_{n+2} = u_n + u_{n+1}$ for all $n \geq 0$. Put

$$U(x) = u_0 + u_1x + u_2x^2 + \cdots = \sum_{n=0}^{\infty} u_n x^n.$$

This formal power series is usually called the *generating function* of u_n , although it is not, strictly speaking, a function. We easily calculate

$$U(x)(1 - x - x^2) = u_0 + (u_1 - u_0)x,$$

since $(u_{n+2} - u_{n+1} - u_n)x^{n+2} = 0$ for all $n \geq 0$.

Hence

$$U(x) = \frac{x}{1 - x - x^2} = \frac{A}{1 - \alpha x} + \frac{B}{1 - \beta x} = A \sum_{n=0}^{\infty} \alpha^n x^n + B \sum_{n=0}^{\infty} \beta^n x^n$$

in partial fractions, where

$$(1 - \alpha x)(1 - \beta x) \equiv 1 - x - x^2,$$

$$A(1 - \beta x) + B(1 - \alpha x) \equiv x.$$

We see from the first identity that $\alpha + \beta = 1$ and $\alpha\beta = -1$, hence

$$\alpha = \frac{1}{2}(1 + \sqrt{5}), \quad \beta = \frac{1}{2}(1 - \sqrt{5}).$$

From the second identity one easily determines

$$A = \frac{1}{\alpha - \beta} = \frac{1}{\sqrt{5}}, \quad B = \frac{1}{\beta - \alpha} = \frac{-1}{\sqrt{5}}.$$

Comparing the two expansions of $U(x)$ above we obtain

$$u_n = A\alpha^n + B\beta^n = \frac{1}{\sqrt{5}} \left(\frac{1 + \sqrt{5}}{2} \right)^n - \frac{1}{\sqrt{5}} \left(\frac{1 - \sqrt{5}}{2} \right)^n.$$

Furthermore, $\beta^n/\sqrt{5}$ is close to 0 and it is not hard to show that u_n is the integer nearest $\alpha^n/\sqrt{5}$.

The number α is called the ‘golden ratio’. It played a role in Euclid’s construction of the regular pentagon. The above formula for u_n was discovered and proved by A. de Moivre (1730). (See Section 9.6 in Stillwell’s *Mathematics and its History*.)

The Fibonacci numbers arose in connection with the following problem from the *Liber abaci*: how many pairs of rabbits will be produced in a year, beginning with a single pair, if in every month each pair bears a new pair which becomes productive from the second month on?

Exercises

1. Solve Gerbert’s problem.
2. Solve Jordanus’s problem.
3. Solve problem (1) on Frederick II’s math contest.
4. Solve problem (2) on Frederick II’s math contest.
5. Starting with the formula $u_n = (\alpha^n - \beta^n)/\sqrt{5}$, show that u_n is the integer nearest $\alpha^n/\sqrt{5}$.
6. Show that u_{n+1}/u_n tends to the golden ratio as n tends to infinity.
7. Show that u_n and u_{n+1} are relatively prime.
8. Prove that the greatest common divisor of any two Fibonacci numbers is also a Fibonacci number.

Mathematics in the Renaissance

Aside from the invention of the Indian numerals, and aside from the work of a few persons of talent, such as Pappus and Fibonacci, no significant advances in mathematics had taken place in the thousand years following Diophantus. In the 15th and 16th centuries there was a sudden spurt of activity, aided by the Chinese invention of printing, which reached Europe in 1450 and which carried mathematics, both pure and applied, beyond the achievements of the ancients. It is hard to overemphasize the importance of printing for the spread of mathematical knowledge. Copying mathematical texts by hand required much time and labour. In ancient times, many texts existed only in a single copy, which would be found in the library of Alexandria. This is why, for about 800 years, all mathematical activity was concentrated in one place. Now such texts were available all over the civilized world and people could learn mathematics even in such outlying places as Bohemia or Scotland. In this chapter, and in the next two chapters, we shall discuss advances in the following areas:

1. mathematical notation,
2. the theory of equations,
3. the invention of logarithms,
4. mechanics and astronomy.

(1) **Mathematical notation**

Johannes Regiomontanus (1436–1476) of Königsberg, then in Germany, gave the first systematic exposition of plane and spherical trigonometry,

using both sines and cosines. In algebra he wrote 'res' for x , and 'census' for the square. Regiomontanus probably died of the plague, but there was a rumour that he was poisoned by the sons of a rival scholar.

Columbus took a copy of Regiomontanus's *Ephemerides* on his fourth voyage, and used its prediction of the lunar eclipse of February 29, 1504 to intimidate some hostile Indians in Jamaica.

Johannes Widman (1462–1500) of Eger, now in the Czech Republic, published a book, *Mercantile Arithmetic*, in 1489, in which the modern symbols $+$ and $-$ appeared for the first time.

Luca Pacioli of Italy (1445–1517) was a Franciscan monk. He used the 'res' and 'census' terminology of Regiomontanus. In 1509 he published the *Divina Proportione* (*Divine (or Golden) Ratio*), a book that was illustrated by none other than Leonardo da Vinci. There is a famous painting of Pacioli by Jacopo de' Barbari, which is now in the Museum of Naples. It shows the friar with his friend Guidobaldo standing in the presence of a dodecahedron. One of the problems solved by Pacioli was the following:

The radius of the inscribed circle of a triangle is 4, and the segments into which one side is divided by the point of contact are 6 and 8. Determine the other sides.

Robert Recorde of England (1510–1558) was the first person to use the symbol $=$ for equality, asserting that 'noe 2 thynges can be more equalle'. Recorde got into a tangle with the Earl of Pembroke and died in jail.

Christoff Rudolff of Germany used $\sqrt{}$ for 'radix' in 1525.

Adam Riese of Bavaria (1492–1559) published arithmetic books that went through more than a hundred editions, and established the use of the signs $+$ and $-$.

Michael Stifel of Germany (1487–1567) was a monk who became an early follower of Luther. He introduced $1A$, $1AAA$, $1AAA$ for A , A^2 and A^3 . He was the first to use negative integers as exponents and had a way of applying mathematics to the Bible which led him to conclude that Pope Leo X was the beast of the *Book of Revelation*, and also to prophesy the end of the world for 18 October, 1533. The peasants of the village where he was pastor believed this prophecy and spent all their money. When the world failed to end, Stifel found himself, not in heaven, but in a jail in Wittenberg.

Thomas Harriot of England (1560–1621) wrote a , aa , aaa for a , a^2 , a^3 and introduced the signs $>$ and $<$ for strict inequality. He went to America with Sir Walter Raleigh and became a tobacco addict. In 1603 Harriot computed the area of a spherical triangle:

Take the sum of all three angles and subtract 180 degrees. Set the remainder as numerator of a fraction with denominator 360

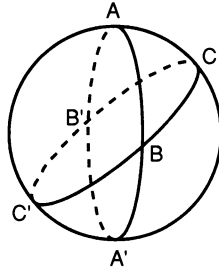


FIGURE 24.1. Area of a spherical triangle

degrees. This fraction tells us how great a portion of the hemisphere is occupied by the triangle.

To prove this, first note that, if a sphere has unit radius, then its surface is 4π or 720° . Hence, the area between two great circles, e.g., between two meridians on the earth, one degree apart, on one hemisphere, is $720^\circ/360 = 2^\circ$. It follows, moreover, that the area between two great circles, on one hemisphere, separated by an angle A is $2A$.

The spherical triangle ABC is bounded by three great circles, as in Figure 24.1, where A' , B' and C' are the antipodes of A , B and C respectively. We note that

$$\begin{aligned}\triangle ABC + \underline{\triangle A'BC} &= 2A, \\ \triangle A'B'C' + \underline{\triangle AB'C'} &= 2B, \\ \underline{\triangle ABC} + \underline{\triangle A'BC'} &= 2C.\end{aligned}$$

The four underlined triangles make up the visible hemisphere (Figure 24.1), namely, 360° ; hence adding the three equations we get

$$\triangle ABC + \triangle A'B'C' + 360^\circ = 2A + 2B + 2C.$$

Now $\triangle A'B'C' = \triangle ABC$, the two being antipodal triangles. Dividing by 2, we obtain

$$\triangle ABC + 180^\circ = A + B + C.$$

Thus the area of the spherical triangle, measured in degrees, is its *spherical excess* $A + B + C - 180^\circ$.

(2) The theory of equations

Ever since the ancient Babylonians, people knew how to find positive, real solutions to any linear or quadratic equation. This they could do arith-

metically or geometrically. Omar Khayyam (1100 A.D.) had developed a method for drawing a line segment whose length was a positive real root of a given cubic polynomial. In 1225, Leonardo of Pisa gave an arithmetical solution to $x^3 + 2x^2 + 10x = 20$. Because he used arithmetic rather than geometry, Leonardo was able to obtain an approximation to the positive root which was accurate to nine decimal places.

The first person to develop anything like a complete method for solving cubic equations — one that could in principle handle negative and imaginary roots as well as positive real roots — was alleged to have been Scipione Ferro of Bologna (1465–1526). He could solve any equation of the form $x^3 + bx = c$, giving the answers to any degree of accuracy. Ferro kept his method a secret — so that he would have an advantage over other mathematicians in mathematical contests — but, just before he died, he passed it on to one Antonio Fiore.

In 1530 Zuanne da Coi sent the following problems to Niccolo Tartaglia (1500–1557):

$$\begin{aligned}x^3 + 3x^2 &= 5, \\x^3 + 6x^2 + 8x &= 1000.\end{aligned}$$

Tartaglia announced that he could solve these equations and was promptly challenged by Fiore to a contest. Each contestant had to deposit a certain amount of money with a notary and propose a number of problems for his rival to solve. Whoever solved more of the problems within 30 days was to get all the money.

Tartaglia, suspecting that Fiore would pose equations of the form $x^3 + bx = c$, quickly worked out a general method for solving such equations. Indeed, Fiore's problems were of this nature, and Tartaglia was able to solve them all. He himself posed equations of the form $x^3 + ax^2 = c$, which he could solve, but which were too difficult for Fiore.

Tartaglia was born in Brescia, Italy. This town was captured by the French in 1512 and most of the inhabitants were massacred. Tartaglia's jaws were split by a soldier's sword and it was thus that he acquired his name, which means 'the stammerer'. He lectured at Verona and Venice, becoming famous through his victory over Fiore. He published a book on ballistics in 1537, in which he correctly stated that a projectile achieves its maximum range when fired at an angle of 45° . However, he gave no proof of this. In 1560 Tartaglia wrote a book on number theory, which contained some amusing puzzles, for example:

Three couples wish to cross a river in a boat that holds only two people. How can this be done if no woman is to be left with a man unless her husband is present?

Three people wish to share the oil in a 24 ounce jar. They have empty measuring jars of capacity 5, 11 and 13 ounces. How can they divide the oil?

The later part of Tartaglia's life was embittered by a quarrel with Girolamo Cardano (1501–1576), another Italian, whose autobiography has been republished by Dover. Cardano was a famous physician in Milan and once travelled to Scotland to cure an Archbishop of asthma. He applied his mathematics to mechanics, gambling and astrology. Indeed, he might be called the discoverer of probability theory. Cardano's oldest son was executed for having poisoned his wife, and Cardano himself was imprisoned, in 1570, for heresy, for having published a horoscope of Jesus Christ. (This was heretical because it suggested that God was subject to the stars.) Cardano was freed only after he recanted. There is a story that he foretold the day of his own death, using astrology, and then felt compelled to commit suicide to make his prediction come true.

Cardano had persuaded Tartaglia to tell him the secret method for solving cubic equations. This Tartaglia did only on condition that Cardano would never reveal it. However, some years later Cardano learned of the prior work by Ferro, and he decided to publish the secret method in his *Ars Magna* (1545). Cardano gave due credit to Tartaglia for the method, but Tartaglia was upset that Cardano had broken his promise to keep it secret. Henceforth, Tartaglia would have no special advantage in mathematical contests!

Annoyed, Tartaglia challenged Cardano to a competition. The latter did not show up, being represented instead by his student, Ferrari (1522–1565). It seems that Ferrari did better than Tartaglia, and Tartaglia lost both prestige and income.

Cardano's *Ars Magna* was the best book on algebra so far. It still used geometry to prove the algebraic identity

$$(a - b)^3 = a^3 - b^3 - 3ab(a - b),$$

and it still shied away from negative numbers, listing the following equations separately:

$$x^3 + px = q, \quad x^3 = px + q, \quad x^3 + px + q = 0, \quad x^3 + q = px.$$

Nonetheless, the *Ars Magna* contained a full explanation of the cubic equation, including a treatment of imaginary numbers.

Ferri came from Bologna in Italy. He extended the work of Tartaglia and Cardano, solving the general fourth degree equation. His solution appears in the *Ars Magna*. Ferrari became rich in the service of the Cardinal Fernando Gonzalo, but ill health forced him to retire to Bologna, in 1565, to teach mathematics. According to W.W. Rouse Ball in *A Short Account*

of the *History of Mathematics*, p. 225, Ferrari was murdered by his sister or by her boyfriend.

Rafael Bombelli of Bologna (1526–1572) published an algebra book in 1572, in which he traced the history of the subject back to Diophantus. He discussed complex radicals at length and showed that the irreducible case of the cubic equation leads to three real roots (see Case 3 Chapter 25). He also pointed out that the ancient Greek problem of trisecting an angle was equivalent to solving a cubic equation. He wrote $\sqrt[n]{}$ for x^n .

Francois Viète, also called Vieta, (1540–1603) was a French lawyer and member of parliament, but his avocation was mathematics. He wrote *In artem analyticem isagoge* in 1591, in which he applied algebra to geometry. (Hitherto people had applied geometry to algebra.) He was challenged by King Henry IV to solve a special equation of the 45th degree, and managed to give the answer in a few minutes, having noticed that the equation was satisfied by the chord of an angle of $360^\circ/45$. He constructed the circles touching three given circles, using only Euclidean geometry, and thus recaptured an ancient construction that was probably contained in a lost book by Apollonius. Viète deciphered a Spanish code for the French. His solutions of cubic and quadratic equations were just like ours.

We close this chapter with a question about astrology. Why did many mathematicians, such as Ptolemy, Cardano and later, Kepler, waste their time and talent on astrology when a little reflection reveals that it seems unlikely to have any truth in it? Astrology is based on the unproven and implausible assumption that there is a correlation between the constellations of the stars at the time of a person's birth and his or her character and ultimate life history. Did Ptolemy, et al., only espouse the practice of astrology because it helped to supplement their incomes, or did they genuinely believe in it, as perhaps a majority of people still do today? Kepler, for one, had a cynical view of astrology, as we shall see.

Exercises

1. Professor Smith learns a secret mathematical technique from Professor Brown only because he solemnly promises Brown that he will never publish it. Later Smith discovers that this technique was long ago published, by Professor Jones, in an obscure little journal that no one ever reads. Is it morally permissible for Brown to publish this technique (giving due credit, of course, to Smith and Jones)? Support your answer with reasons.

2. Solve Pacioli's problem, given above.
3. Solve Tartaglia's puzzle about the three couples.
4. Solve Tartaglia's puzzle about the oil.
5. Let $a_1 = 1/\sqrt{2}$, $a_{n+1} = \sqrt{\frac{1}{2} + \frac{1}{2}a_n}$. Viète proved that

$$2/\pi = a_1 a_2 \cdots a_n \cdots$$

Prove this formula. (Hint: first show

$$(\sin x)/(2^n \sin(x/2^n)) = (\cos x/2)(\cos x/2^2)(\cos x/2^3) \cdots (\cos x/2^n),$$

whence

$$(\sin x)/x = (\cos x/2)(\cos x/2^2)(\cos x/2^3) \cdots$$

and, finally, let $x = \pi/2$.)

The Cubic and Quartic Equations

Cardano was the first person to use imaginary numbers in print. In this chapter, we shall use imaginary numbers to present what is essentially Cardano's solution to the cubic equation. We shall also give Ludovico Ferrari's solution to the fourth degree polynomial equation.

Recall that $i = \sqrt{-1}$. Recall that, if a, a', b and b' are real numbers and $a + bi = a' + b'i$, then $a = a'$ and $b = b'$ (lest $i = (a - a')/(b - b')$, a real number). Let $\omega = \frac{-1}{2} + \frac{1}{2}\sqrt{3}i$. A quick calculation shows that $\omega^2 = \frac{-1}{2} - \frac{1}{2}\sqrt{3}i$, and hence $\omega^3 = 1$.

Lemma 25.1. *Any complex number $x + iy$ can be written in the polar form $r(\cos A + i \sin A)$, where r is a non-negative real number and A is a real number.*

Lemma 25.2. $(r(\cos A + i \sin A))^3 = r^3(\cos 3A + i \sin 3A)$.

Proof:

$$\cos 3A = \cos^3 A - 3 \cos A \sin^2 A,$$

$$\sin 3A = 3 \cos^2 A \sin A - \sin^3 A.$$

Lemma 25.3. *The equation $z^3 = 1$ has exactly three complex solutions: $1, \omega, \omega^2$.*

Lemma 25.4. *If b is any given complex number $\neq 0$, the equation $z^3 = b$ has exactly three complex solutions. If z_1 is one solution, the others are $z_1\omega$ and $z_1\omega^2$.*

Lemma 25.5. *If y and p are any complex numbers, there are complex numbers u and v such that $u + v = y$ and $uv = p$. (This is the old Babylonian*

problem.)

The proofs of these Lemmas are left to the reader.

The general cubic equation (after division by the leading coefficient) has the form

$$x^3 + ax^2 + bx + c = 0.$$

(We assume here that a, b , and c are real.) Putting $x = y + k$, this equation becomes

$$y^3 + (3k + a)y^2 + (\dots)y + (\dots) = 0.$$

We choose $k = -a/3$, so that the square term disappears. (Tartaglia was the discoverer of this trick.) The resulting equation has the form

$$y^3 - 3py - 2q = 0$$

with p and q real and where the numbers -3 and -2 are introduced for convenience only. It is this *reduced* equation which we now want to solve.

If we put $y = u + v$, then the reduced equation becomes

$$u^3 + v^3 + 3(uv - p)(u + v) - 2q = 0.$$

To make this as simple as possible, we choose $v = p/u$ (see Lemma 25.5). The equation then becomes

$$u^3 + v^3 = 2q.$$

Since, however, $u^3v^3 = p^3$, we are back with the Babylonian problem of seeking two numbers with given sum and product. Clearly u^3 and v^3 are solutions of

$$t^2 - 2qt + p^3 = 0.$$

Therefore we have, say,

$$u^3 = q + \sqrt{q^2 - p^3},$$

$$v^3 = q - \sqrt{q^2 - p^3}.$$

Thus $y = u + v = u + p/u$, where u is a cube root of

$$q + \sqrt{q^2 - p^3}.$$

If one of these cube roots is u_1 , then the others are $u_1\omega$ and $u_1\omega^2$. Let $v_1 = p/u_1$. Then $u_1\omega + p/(u_1\omega) = u_1\omega + v_1\omega^2$ and $u_1\omega^2 + p/(u_1\omega^2) = u_1\omega^2 + v_1\omega$. Hence the reduced cubic has the following three solutions only:

$$u_1 + v_1, \quad u_1\omega + v_1\omega^2, \quad u_1\omega^2 + v_1\omega.$$

To discuss these solutions in detail, we consider three cases.

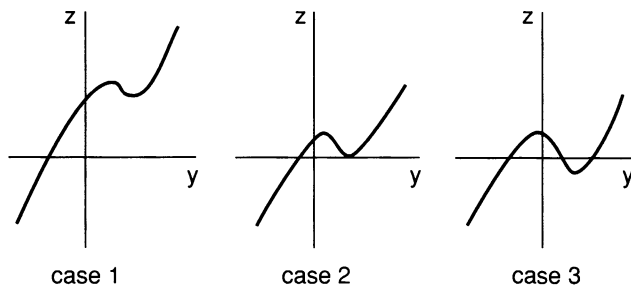


FIGURE 25.1. Cubic polynomials

Case 1. $q^2 - p^3 > 0$. Say $q^2 - p^3 = r^2$, where r is real. Let u_1 be the real cube root of $q + r$, so that v_1 is the real cube root of $q - r$ (the product $u_1 v_1$ must be real). Then the solutions are the real number $u_1 + v_1$ and the imaginary ‘conjugates’ $u_1 \omega + v_1 \omega^2$ and $u_1 \omega^2 + v_1 \omega$. (Note that ω^2 is the *conjugate* of ω , that is, it may be obtained from ω by replacing i with $-i$.)

Case 2. $q^2 - p^3 = 0$. Then we can take u_1 as the real cube root of q , and v_1 the same. The equation then has only two distinct roots, namely, $2u_1$ and $-u_1$ (since $\omega + \omega^2 = -1$), but we say that the latter root occurs twice.

Case 3. $q^2 - p^3 < 0$. Say $q^2 - p^3 = -r^2$, where r is a positive real number. Then $u^3 = q + ir$ and $v^3 = q - ir$. Let $u_1 = a + bi$ with a and b real. Then $(a^2 + b^2)^3 = q^2 + r^2 = p^3$, so that $a^2 + b^2 = p$ and hence $v_1 = p/u_1 = a - bi$. We calculate

$$u_1 \omega + v_1 \omega^2 = a(\omega + \omega^2) + bi(\omega - \omega^2) = -a - b\sqrt{3}.$$

Similarly, $u_1 \omega^2 + v_1 \omega = -a + b\sqrt{3}$. It is not hard to show that $a = \sqrt{p} \cos((\arctan r/q)/3)$ and $b = \sqrt{p} \sin((\arctan r/q)/3)$. (Note that, in Case 3, $p > 0$.)

If $z = y^3 - 3py - 2q$, the three cases are illustrated by the three graphs in Figure 25.1. However, it should be borne in mind that these graphs were not available in the Renaissance, as analytic geometry had not yet been invented.

The general quartic equation has the form

$$x^4 + ax^3 + bx^2 + cx + d = 0,$$

which we may write

$$x^4 + ax^3 = -bx^2 - cx - d.$$

Adding $a^2x^2/4$ to both sides of this equation, we obtain

$$\left(x^2 + \frac{1}{2}ax\right)^2 = \left(\frac{1}{4}a^2 - b\right)x^2 - cx - d.$$

In an attempt to get a perfect square on both sides of the equation, we add $t(x^2 + \frac{1}{2}ax) + \frac{1}{4}t^2$ to both sides:

$$\left(x^2 + \frac{1}{2}ax + \frac{1}{2}t\right)^2 = \left(\frac{1}{4}a^2 - b + t\right)x^2 + \left(-c + \frac{1}{2}at\right)x - d + \frac{1}{4}t^2.$$

Now $Ax^2 + Bx + C$ is a perfect square when $B^2 - 4AC = 0$. In fact, if $A \neq 0$, we can then write

$$Ax^2 + Bx + C = (\sqrt{A}x + B/(2\sqrt{A}))^2.$$

So, to get a square, it will suffice to pick t so that

$$-d + \frac{1}{4}t^2 = \frac{\left(-c + \frac{1}{2}at\right)^2}{4\left(\frac{1}{4}a^2 - b + t\right)},$$

that is,

$$t^3 - bt^2 + (ac - 4d)t + 4bd - a^2d - c^2 = 0.$$

But this is a cubic equation! In practice, this associated cubic equation can often be solved by trial and error. We only need one value of t .

We see that if x is such that $x^4 + ax^3 + bx^2 + cx + d = 0$, then there is a t , which we can determine by finding a real root of the above cubic equation, such that

$$\left(x^2 + \frac{1}{2}ax + \frac{1}{2}t\right)^2 = (\sqrt{A}x + B/2\sqrt{A})^2,$$

with $A = \frac{1}{4}a^2 - b + t$, $B = -c + \frac{1}{2}at$ and $C = -d + \frac{1}{4}t^2$. This gives us

$$x^2 + \frac{1}{2}ax + \frac{1}{2}t = \pm(\sqrt{A}x + B/2\sqrt{A}),$$

which is a quadratic equation in x .

For several centuries people tried to find similar methods for solving equations of degree greater than 4. They failed. It was only in the 19th century that it was shown by Ruffini, Abel and Galois that the general equation of degree 5 or more cannot be solved by 'radicals' (e.g. fifth roots). Of course, special cases can be solved by radicals, for example $x^5 + x = 34$.

Exercises

1. Prove all the Lemmas in this chapter.
2. Solve the following equations, obtaining exact answers. Do not use decimal approximations. Simplify answers.
 - (a) $x^3 + x^2 - 2 = 0$,
 - (b) $x^3 + 9x - 2 = 0$,
 - (c) $y^3 - 3y + 1 = 0$,
 - (d) $y^3 - 7y - 7 = 0$,
 - (e) $x^3 + 2x^2 + 10x = 20$,
 - (f) $x^3 + 3x^2 = 5$,
 - (g) $x^3 + 6x^2 + 8x = 1000$,
 - (h) $x^4 + x^3 - 6x^2 - x + 1 = 0$,
 - (i) $x^3 = 15x + 4$.
3. Why cannot a fourth degree equation have five or more distinct complex roots?

Renaissance Mathematics Continued

(3) The invention of logarithms

John Napier (1550–1617), a Scottish aristocrat, spent 20 years of his life on the construction of logarithms. Like Stifel, he was interested in proving that the Pope was the Antichrist (the opponent of Jesus, who is supposed to appear just prior to the latter’s prophesied second coming). Napier was convinced that the world would end before 1700.

Napier describes his technique for calculating logs in his *Mirifici Logarithmorum Canonis Constructio*, which was published in 1619, two years after the author’s death.

This *Construction of the Wonderful Canon of Logarithms* is tricky to read, because what Napier calls the ‘logarithm’ of x is actually what we call $10^7 \log_{1/e}(x/10^7)$. The book is also a bit tedious because Napier pays minute attention to error bounds. However, if we simplify and modernize Napier’s presentation somewhat — as we shall do below — we obtain a very lucid piece of mathematics.

Napier’s basic ideas are these. Suppose there is a particle on the negative half of the real number line, moving towards the origin at a speed proportional to its distance from the origin. At time 0, it is at -1 . For any distance d , there is a time when the particle is that distance from the origin. Call this time the ‘logarithm’ of d .

Assume that the constant of proportionality is 1. Then $dx/dt = -x$ and $x(0) = -1$. Hence $x(t) = -e^{-t}$, and so $t = -\log_e d = \log_{1/e} d$. In this section ‘ $\log x$ ’ shall mean $\log_{1/e} x$. Note that $\log_{1/e} d$ decreases as d increases from $1/2$ to 1.

Napier used the following three ideas to calculate his tables.

1. When z is very small, $\log(1 - z) \approx z$.
2. The natural number powers of $(1 - 10^{-m})$ — where m is a natural number — are easy to calculate. For, if $s = 1 - 10^{-m}$, then $s^2 = s(1 - 10^{-m}) = s - s/10^m$. Similarly, $s^3 = s^2 - s^2/10^m$. It is always just a matter of shifting the decimal point m places and subtracting. Thus Napier has

$$\begin{aligned}
 1 - 10^{-5} &= 0.999\,990\,000\,000\,000, \\
 (1 - 10^{-5})^2 &= 0.999\,990\,000\,000\,000 \\
 &\quad - .000\,009\,999\,900\,000 \\
 &= 0.999\,980\,000\,100\,000, \\
 (1 - 10^{-5})^3 &= 0.999\,980\,000\,100\,000 \\
 &\quad - .000\,009\,999\,800\,001 \\
 &= 0.999\,970\,000\,299\,999.
 \end{aligned}$$

3. Near 1, the log is smooth, and linear interpolations give excellent approximations to it.

Near -1 , Napier's particle is moving at about 1 unit/second, and we can assume that $\log d = 1 - d$ for $d = 1 - 10^{-5} = 0.99999$. As in (2) above, Napier calculates d, d^2, \dots, d^{50} . He finds that $d^{50} = 0.9995001225$. Hence $\log 0.9995001225 = 50 \times \log d = 50 \times 0.99999$. With linear interpolation, Napier can thus obtain a very accurate value for $\log u$, where $u = 0.9995$.

Next, using ideas similar to those in (2), Napier quickly and accurately calculates u^2, u^3, \dots, u^{20} . He obtains $u^{20} = 0.990047358$. Using interpolation, he then gets a very accurate value for $\log w$, where $w = 0.99$. (Your pocket calculator will not be more accurate.)

For $a = 1, 2, \dots, 20$ and $b = 0, 1, \dots, 68$, Napier calculates $u^a w^b$. This gives him 1380 points between $u = 0.9995$ and $u^{20} w^{68} = 0.499860940$ for calculating logarithms in that range.

For example, $u^{19} w^{68} = 0.500110996 > 1/2 > 0.499860940 = u^{20} w^{68}$. If

$$k = \frac{\log(u^{19} w^{68}) - \log(u^{20} w^{68})}{u^{19} w^{68} - u^{20} w^{68}},$$

then k is the slope of the line joining two given points. Thus, by linear interpolation we have

$$\begin{aligned}
 \log 1/2 &\approx \log(u^{20} w^{68}) + (\tfrac{1}{2} - u^{20} w^{68})k \\
 &\approx 20 \log u + 68 \log w - 0.000278 \\
 &\approx 0.693147.
 \end{aligned}$$

This value of $\log 1/2$ is accurate to six decimal places.

Of course, it is now easy to calculate other logs. For example, if $0 < t < 1/2$, we can calculate $\log t$ by finding an integer m such that $1/2 < 2^m t < 1$. Given such an m , we have $\log t = \log(2^m t) + m \log 1/2$. The \log of 2 (to our base $1/e$) is just $-\log 1/2$.

Henry Briggs (1561–1631) travelled to Edinburgh in 1615 to discuss logarithms with Napier. They agreed that there were many advantages to having logs to the base 10. In 1617 Napier died, but Briggs continued the work, publishing tables for logs to the base 10 in 1624. It was Briggs who introduced the terms ‘mantissa’ and ‘characteristic’. The practical advantage of base 10 is of course that the logarithm of numbers such as 173, 17.3, 1.73, 0.173 and 0.0173 all have the same mantissa, which can be found in the table, while their characteristics 2, 1, 0, -1 , and -2 are seen by inspection. Briggs persuaded many of his contemporaries, including Kepler, of the importance of logarithms.

(4) **Mechanics and astronomy**

Nicholas Copernicus (1473–1543) was of Polish origin. He conjectured that the earth and the other planets move around the sun. At any rate, he said that this assumption gives a simpler explanation of what is going on. He did not assert that the planets move in circles with the sun at the center. He could only describe the orbits as epicycles, as did Ptolemy, whose observational data he used.

Simon Stevin (1548–1620) lived in the Netherlands (which included what is now Belgium). He was one of the earliest expositors of the theory of decimal fractions, which he compared to an unknown island ‘having beautiful fruits, pleasant plains, precious minerals’ (see page 21 in Smith’s *A Source Book in Mathematics*). Stevin is perhaps best known for his *Statistics and Hydrostatics*, published in 1586. In this book he discusses the triangle of forces, resolution of forces, stable and unstable equilibrium and pressure. He was the first to advance beyond Archimedes in these subjects. He also wrote on algebra and geometry.

Galileo Galilei (1564–1642) is often regarded as the father of modern physics. There is a story that, instead of paying attention to the church service, he used his pulse to time the oscillations of a lamp swinging from the church roof. He discovered that the period of oscillation does not depend on its amplitude. This discovery was to be exploited for constructing pendulum clocks. Galileo disrupted his medical studies to devote himself to mathematics. After publishing a book on hydrostatics and centers of gravity, he was appointed professor. Everyone has heard the story about his experimenting with falling bodies on the leaning tower of Pisa. Whatever the truth of this story, he found that a falling body undergoes a uniform ac-

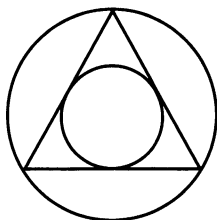


FIGURE 26.1. Kepler's solar system

celeration which is independent of its weight. He checked these conclusions by studying bodies sliding on inclined planes.

Galileo is most famous as an astronomer. He was the first to use the recently invented telescope to study the heavens and was rewarded by discovering, among other things, the moons of Jupiter (in 1610). These seemed to confirm the Copernican hypothesis which Galileo advocated. It is well-known that, in 1633, the Inquisition forced him to recant the Copernican view.

Johann Kepler (1571–1630) studied in Tübingen (Germany) and became a professor in Grätz (Austria). His early espousal of the Copernican system led him to seek an explanation of the distances of the various planets from the sun. His first idea was to construct an equilateral triangle with its vertices on the orbit of Saturn — it was assumed that this orbit was circular — and then to inscribe another circle in this triangle (Figure 26.1).

This second circle was supposed to be the orbit of Jupiter, and, indeed, this construction gave more or less the correct ratio for the distances of Saturn and Jupiter from the sun. Kepler tried to extend this idea by constructing a square with its vertices on the orbit of Jupiter, and inscribing a circle in this square for the orbit of Mars. Here, however, the observed facts did not confirm to the theoretical model.

Undiscouraged, Kepler then replaced the circles by spheres and the regular polygons by regular polyhedra. Remarkably, there were five intervals between the known planetary orbits to account for and five Platonic solids to explain them. Using the observational data of Ptolemy, Kepler made his model fit the facts more or less, and he published his findings in his *Cosmic Mystery*, a book which owed much to Pythagorean number mysticism. Today we do not attribute any significance to these speculations.

Kepler knew that, in order to check his theories, he had to use more recent observational data than those recorded by Ptolemy. It so happened that very accurate observations of celestial phenomena were being made by the Danish astronomer Tycho Brahe (1546–1601), who was, however, reluctant to part with his data, hoping to save them for his own theory, which was a modification of Ptolemy's.

The capital of the Holy Roman Empire was at that time in Prague. Tycho Brahe was called to Prague as ‘imperial mathematicus’ and Kepler managed to become his assistant in 1599. Brahe put Kepler to work on the orbit of Mars, which deviated from a circle more than any other planetary orbit did. Two years later, Brahe died, and Kepler succeeded him as imperial mathematicus. His job description included casting horoscopes for the emperor. Kepler expressed his views on astrology in *Tertius Interveniens*, asserting that it is all very well for philosophers to criticize the ‘daughter’ of astronomy, without realizing that ‘the daughter must support the mother by her charms’. He pointed out that an astronomer could not make a living, unless he encouraged people in the belief that they could learn the future from the stars.

With the help of Brahe’s magnificent data, Kepler was able to formulate his three laws of planetary motion:

1. Each planet describes an ellipse with the sun at one focus.
2. The line joining the sun to the planet sweeps out equal areas (bounded by the ellipse) in equal times.
3. The square of the period of revolution of each planet (its ‘year’) is proportional to the cube of the major axis of its orbit.

Kepler’s third law came in 1619, about ten years after the first two. These three laws were later to confirm Newton’s theory of universal gravitation.

After the job in Prague petered out, owing to political events and a shortage of funds in the imperial treasury, Kepler supported himself by casting horoscopes. He used mathematics even in his private life, making a careful calculation to decide which eligible woman to choose as his second wife.

There was to be one occasion when Kepler’s connection to the court in Prague proved useful. When jealous neighbours had accused his mother of witchcraft, a very serious accusation in those days, he managed to get the case dismissed.

Exercises

1. How should Napier have had his particle moving to get logs to the base e ?
2. Suppose we have perfectly accurate values for $\log(u^a v^b)$ and $\log(u^{a-1} v^b)$. Suppose $u^a v^b < x < u^{a-1} v^b$, and suppose we calculate the value of $\log x$ from the values of the two given logs, using linear interpolation. What is the maximum error possible?

3. You have been stranded on a desert island without your calculator. Being very bored, you calculate 2^{1000} , doubling 999 times but keeping only five significant figures (leaving a trail of zeros at the right of the calculation). You find that $2^{1000} = 1.0715... \times 10^{301}$. By what easy way can you now find $\log_{10} 2$, accurate to three decimal places?
4. Briggs calculated logs using square roots. He found $10^{1/2}$, $10^{1/4}$, $10^{1/8}$ etc. He then found numbers such as $10^{3/8} = 10^{1/4}10^{1/8} = 2.37...$. This gave him the $\log_{10} 2.37...$. Use Briggs's method, together with linear interpolation, to get a value for $\log_{10} 2.37...$.
5. According to Kepler's initial model, what is the ratio of the radius of Saturn's orbit to the radius of Jupiter's orbit?

The Seventeenth Century in France

The 17th century saw a blossoming of mathematical activity in France. Some of the important mathematicians were:

- Marin Mersenne (1588–1648),
- Gerard Desargues (1591–1661),
- René Descartes (1596–1650),
- Pierre de Fermat (1601–1665),
- Blaise Pascal (1623–1662).

Mersenne was a friar belonging to the Minim order. In the chapter on perfect numbers, we learned about primes of the form $2^n - 1$, which are named after him. His main importance lies in the fact that he corresponded with all the other French mathematicians, keeping them in touch with each other's ideas.

Desargues was the discoverer of projective geometry, the part of the geometry which deals entirely with incidence and ignores distance. Parallel lines are presumed to meet at a point 'at infinity'. In retrospect, it appears that the theorem of Pappus was also a theorem in projective geometry. Desargues's famous theorem is this:

If two triangles, in the same plane or not, are so situated that lines joining pairs of corresponding vertices are concurrent (if

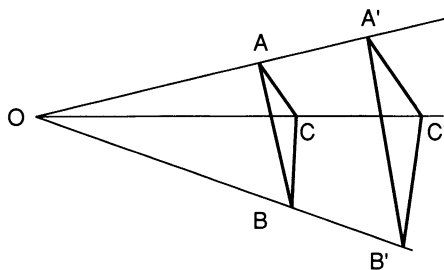


FIGURE 27.1. Desargues's Theorem

only at infinity), then the points of intersection of pairs of corresponding sides are collinear (if only on the 'line at infinity'), and conversely.

The 'parallel case' for two triangles in the same plane is pictured in Figure 27.1. According to the theorem, if the lines joining corresponding vertices meet in a single point O and if AB is parallel to $A'B'$ and AC is parallel to $A'C'$, then BC is parallel to $B'C'$.

The theorem of Desargues can be proved, in Euclidean plane geometry, using only axioms I to V of Chapter 16. This is shown by David Hilbert (1862–1943) in his *Foundations of Geometry*.

Hilbert also shows that, in the absence of the Axiom of Archimedes (mentioned in Chapters 17 and 19), the Theorem of Desargues can be used for defining multiplication within Euclidean geometry, and for proving, within that geometry, that multiplication is associative. (Hilbert also shows that the Theorem of Pappus does not depend on the Axiom of Archimedes. The Theorem of Pappus is the key ingredient in any proof of the commutativity of multiplication which does not rest on the Axiom of Archimedes.)

For a simple proof of the 'parallel case' of Desargues' Theorem, the reader may wish to consult, e.g., Ewald [1971].

Descartes was educated by the Jesuits, who kindly permitted him to spend the mornings in bed, because of his delicate health. Throughout his life, he did much of his intellectual work in bed. Being a man of 'good' family, he was supposed to choose the church or the army for his career, and he chose the latter. He first served in the army of Maurice of Orange, but transferred to that of the Duke of Bavaria when the Thirty Years War broke out. Even on military campaigns, he spent a good part of his time in bed, thinking about mathematics and philosophy. He resigned his commission in 1621, and travelled for five years, studying mathematics. In the end, he settled in Holland, where he spent twenty years in full time intellectual pursuits. In 1649 he went to Sweden at the invitation of Queen Christina. The vigorous young queen insisted on having mathematics lessons at five in the morning, and within two months the aging Descartes, who would

have preferred to sleep in, caught pneumonia and died.

Descartes's first book, *Le Monde*, advanced the Copernican model of the solar system. Just as he was completing it, the Inquisition condemned the Copernican views expressed by Galileo. Descartes decided not to publish his book, writing sadly to Mersenne:

This has so strongly affected me that I have almost resolved to burn all my manuscript, or at least show it to no one. But on no account will I publish anything that contains a word that might displease the Church.

In 1637 Descartes published his famous *Discours de la Méthode pour bien conduire sa Rasion et chercher la Verité dans les Sciences*. This 'discours' has three important appendices. Appendix I, *La Dioptrique* treats optics and the laws of refraction. Appendix II, *Les Météores*, deals with atmospheric phenomena; in particular, it offers an explanation of the shape of the rainbow. For us, the most important is Appendix III, *La Géométrie*. This is divided into three books and sets forth the principles of analytic geometry. In it we find the usual formulas for the conics, and also Descartes's 'rule of signs' (stated without proof). Descartes was thus the creator of analytic geometry, although Fermat must share equal credit.

The revolutionary idea of analytic geometry was this: points in the Euclidean plane could be represented by pairs of real numbers and, consequently, straight lines and conic sections could be described by sets of pairs (x, y) satisfying equations of the form

$$Ax + By + C = 0$$

and

$$Ax^2 + Bxy + Cy^2 + Dx + Ey + F = 0,$$

respectively. Thus geometry was reduced to algebra! It was now possible to solve some of the problems which the Greeks had left open, for instance, the problem of doubling the cube, although the actual solution had to wait a couple of centuries more.

Descartes's rule of signs asserts that the number of positive real roots of a polynomial equation $f(x) = 0$ with real coefficients is $v - 2k$, where v is the number of *variations in sign* and k is some natural number. To calculate v one writes $f(x)$ in descending powers of x , omitting all terms with zero coefficients. Then v is the number of times we change sign as we go from left to right. For example, the equation

$$(1) \quad x^6 - 3x^5 - x^4 + 2x - 5 = 0$$

has three variations in sign, so we can tell that it has either three positive roots or only one. Sometimes we can be certain; for example, putting $x = -y$ in equation (1), we get

$$(2) \quad y^6 + 3y^5 - y^4 - 2y - 5 = 0,$$

which has one variation in sign, hence exactly one positive root. It follows that equation (1) has exactly one negative root. Descartes's rule of signs is a special case of what logicians call *elimination of quantifiers*. Descartes wrote two other books, the *Meditations* (1641) and *Principia Philosophicae* (1644). The latter deals with physical science and proposes the theory of vortices to explain planetary motion. This theory was later refuted by Newton.

The *Meditations* contains a 'geometrical proof' of the existence of God:

existence can no more be separated from the essence of God than the idea of a mountain from that of a valley, or the equality of its three angles to two right angles, from the essence of a triangle (Meditation V).

In other words, God exists because existence is just one of the defining properties of God.

Fermat was a councillor for the parliament of Toulouse, and only did mathematics in his spare time. He published almost nothing during his lifetime; his contributions to mathematics are contained in his correspondence (e.g., with Mersenne) and in the papers that were found after his death. Nonetheless, he is considered to be the greatest amateur mathematician of all times.

Fermat introduced his version of analytic geometry in his *Ad Locos Planos et Solidos Isagoge*. He also collaborated on probability theory with Pascal. In addition, Fermat studied tangents to curves, maxima and minima, and areas under curves, coming very close to discovering calculus. Fermat's methods were influenced by those of his contemporaries, Cavalieri and Wallis, whom we shall meet in the next Chapter.

Today Fermat is best known for his contributions to the theory of numbers. Diophantus knew that a prime number of the form $4n - 1$ is never the sum of two squares. Fermat went further and proved that every prime of the form $4n + 1$ can be written as the sum of two squares in exactly one way. He also noted that every odd prime p can be written as the difference of two squares in one and only one way.

The last statement is easy to prove. Indeed, $p = (\frac{1}{2}(p+1))^2 - (\frac{1}{2}(p-1))^2$ where $\frac{1}{2}(p+1)$ and $\frac{1}{2}(p-1)$ are both integers, since p is odd. On the other hand, if $p = x^2 - y^2 = (x+y)(x-y)$, then $x+y = p$ and $x-y = 1$ (since p , being prime, has no factors except 1 and p). This gives $x = \frac{1}{2}(p+1)$ and $y = \frac{1}{2}(p-1)$ as before.

Fermat's so-called 'Little Theorem' asserts that, if p is a prime and a is an integer which is not a multiple of p , then $a^{p-1} - 1$ is a multiple of p . There are many proofs of this theorem. One of them depends on the

Binomial Theorem (known to the Chinese centuries before Fermat):

$$(x+1)^p = x^p + \binom{p}{1}x^{p-1} + \binom{p}{2}x^{p-2} + \cdots + \binom{p}{p-1}x + 1,$$

where

$$\binom{p}{k} = \frac{p(p-1)\cdots(p-k+1)}{k(k-1)\cdots 1}$$

is an integer. If $0 < k < p$, then p is not a factor of $k(k-1)\cdots 1$. Since p is a factor of

$$\binom{p}{k}k(k-1)\cdots 1 = p(p-1)\cdots(p-k+1)$$

it follows that p is a factor of $\binom{p}{k}$ when $0 < k < p$. Hence

$$(x+1)^p = x^p + 1 + mp$$

for some integer m . If $x = 1$, we obtain $2^p = 2 + mp$, so that p is a factor of $2^{p-1} - 1$. If $x = 2$, we get $3^p = 2^p + 1 + m'p = 2 + mp + 1 + m'p = 3 + (m + m')p$, so that p is a factor of $3^{p-1} - 1$. Continuing in the same way (using mathematical induction), we see that, for all a , $a^p = a + np$ for some integer n . Thus p divides $a^p - a = a(a^{p-1} - 1)$ and the theorem follows (since p does not divide a).

More famous still is Fermat's 'Last Theorem'. This was proved in 1994 by Andrew Wiles.

What is Fermat's 'Last Theorem'? In reading Bachet's translation (from Greek into Latin) of the work of Diophantus, Fermat came across the equation $x^2 + y^2 = z^2$ with its solutions (e.g., $x = 3, y = 4$ and $z = 5$). Fermat wrote in the margin of the book that he had been able to prove that the equation

$$x^n + y^n = z^n$$

has no positive integer solutions for $n > 2$, but that the margin was too small for him to write the proof there. (Of course, $0^3 + 7^3 = 7^3$, but here one of the integers is 0; this is called a 'trivial' solution.)

It is quite possible that Fermat had a proof of the nonexistence of the positive integer solutions for $n = 3$. We still have his proof of the 'Last Theorem' for the special case $n = 4$. However, it is unlikely that he ever had a proof for the complete theorem.

Legendre disposed of the case when $n = 5$ in 1823, and Dirichlet handled the case when $n = 14$ in 1832. In 1849 Kummer made a big step forward and was able to vindicate Fermat's statement for all $n < 100$ except 37, 59 and 67. Before 1994, thanks to the help of the computer, we knew that Fermat was right for all $n < 10^8$ or so, but a proof of the general theorem still escaped us. In 1994, however, Wiles gave a complete proof, for all n .

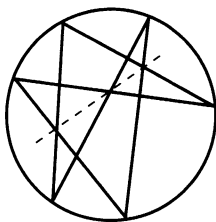


FIGURE 27.2. Pascal's Theorem

Over the centuries, research into Fermat's Last Theorem proved very fruitful. As one example, Kummer's research led to the discovery of 'ideals' and their unique factorization. As another example, Faltings's research led to advances in algebraic geometry — this as recently as 1983. In fact, Faltings came close to proving Fermat's Last Theorem, by showing that, for any $n > 2$, the equation $x^n + y^n = 1$ has at most a finite number of rational solutions.

Pascal was educated at home and forbidden to study mathematics. By the age of 12, he had rediscovered many of Euclid's theorems, so his father relented and gave him a copy of the *Elements*. At 14 Pascal began attending meetings of a group of mathematicians which included Mersenne. At 16 he wrote an essay on conic sections (with an account of the 'mystical hexagram') and, at 18, he constructed a calculating machine, one of the first computers.

At 27, Pascal abandoned mathematics to devote himself wholly to the philosophy of religion and the worship of God. At least one historian of mathematics (E.T. Bell) has taken this as a sure indication of insanity, but it is unlikely that a madman could have produced the elegant French prose or the brilliant philosophical analyses that we find in the *Pensées*, which Pascal wrote during this period. In his later life, Pascal returned to mathematics for a couple of brief periods. Pascal died at age 39. His last feat was the creation of a public transportation system, with the profits going to help the poor.

What did Pascal do as a mathematician? In 1639 (at age 16) he used some ideas of Desargues to obtain 'Pascal's Theorem' (with the 'mystic hexagram'):

Theorem 27.1. *If a hexagon is inscribed in a conic, the points of intersection of opposite sides, assumed to be nonparallel, lie in a straight line (Figure 27.2).*

Actually, this is a theorem in projective geometry, and it suffices to prove it for the case when the conic is a circle, since other conics can be obtained

from the circle by projection. If one views a pair of straight lines as a degenerate conic, one obtains the Theorem of Pappus as a special case of the Theorem of Pascal. A simple proof for Pascal's Theorem in the case of the circle is found in Coxeter and Greitzer's *Geometry Revisited*.

In 1653, Pascal rediscovered what we call 'Pascal's triangle':

$$\begin{array}{ccccccc}
 & & & & 1 & & & & \\
 & & & & 1 & & 1 & & \\
 & & & 1 & & 2 & & 1 & \\
 & & 1 & & 3 & & 3 & & 1 \\
 1 & & 4 & & 6 & & 4 & & 1 \\
 & \cdot & \cdot & \cdot & \cdot & \cdot & & &
 \end{array}$$

As we noted, this goes back to the Chinese, but Pascal was the first to give clear and complete demonstrations of its basic properties (making the first explicit use of mathematical induction). He showed that the k th entry in the n th row is $\binom{n-1}{k-1}$, this being the number of ways of choosing $k-1$ out of $n-1$ things. (Today we prefer to call $\binom{n}{k}$ the number of k -element subsets of an n -element set.) Pascal showed that

$$\binom{n-1}{k-1} = \frac{(n-1)(n-2)\dots(n-k+1)}{(k-1)!}$$

and he showed that

$$\binom{n}{k} = \binom{n-1}{k} + \binom{n-1}{k-1}.$$

(This latter equation follows from the observation that, to choose k out of n things, we may set one of the n things aside and consider two cases: in the first case, the special thing is included in the choice, and hence there are $\binom{n-1}{k-1}$ possibilities; in the second case the special thing is excluded, leaving $\binom{n-1}{k}$ possibilities.) Finally, Pascal proved that $\binom{n}{k}$ is the coefficient of $x^k y^{n-k}$ in the binomial expansion of $(x+y)^n$.

With Fermat, Pascal developed the theory of probability, including the concept of mathematical expectation. Pascal uses probability in the *Pensées* as part of a proof that it is wiser to believe in God (Pascal, *Oeuvres Complètes*, p. 1212). This argument, called 'Pascal's wager', goes as follows: Even if the probability that a god exists is very small (but positive), if he rewards those who believe in him with eternal happiness, assumed to be of infinite value, a rational human being ought to believe in this god. (Did Pascal consider the possibility that God might punish those who adopt their beliefs in expectation of personal gain?)

Exercises

1. Assuming the usual theory of similar triangles (which Euclid based on the Axiom of Archimedes), prove the parallel case of Desargues's Theorem, and its converse. (The converse is: if $AB, A'B'$ and $AC, A'C'$ and $BC, B'C'$ are three pairs of parallels then AA', BB' and CC' are concurrent.)

2. Descartes showed that, roughly speaking, a curve is a conic if and only if its equation has the form

$$Ax^2 + Bxy + Cy^2 + Dx + Ey + F = 0.$$

Give examples of equations of this form which represent the following curves:

- (a) a circle,
 - (b) two parallel straight lines,
 - (c) a point,
 - (d) the whole plane.
3. Why cannot an integer of the form $4n - 1$ be the sum of two integer squares?
 4. What is the smallest positive integer which quadruples when its final (base ten) digit is shifted to the front? (For example, it is not 125 since 512 is not 4 times 125.)
 5. Prove that, if $x^p + y^p = z^p$ has one solution in positive integers, then it has infinitely many such solutions.
 6. Show that, if $x^p + y^p = z^p$, with p an odd prime, then p is a factor of $x + y - z$.
 7. Give a proof of the Binomial Theorem.
 8. Find two distinct proofs to show that the sum of the entries in any row in Pascal's triangle is a power of 2.
 9. Show that the alternating (adding and subtracting) sum of the entries in any row of Pascal's triangle is 0.
 10. Using mathematical induction, show that the $(n + 1)$ th Fibonacci number is

$$\binom{n}{0} + \binom{n-1}{1} + \binom{n-2}{2} + \cdots.$$
 11. Project: Prove Pascal's mystic hexagram theorem for the case of the circle. Note that some of the points of intersection may be outside the circle. (Hint: use the Theorem of Menelaus.)

The Seventeenth Century Continued

Mathematics in the 17th century was by no means confined to France. Elsewhere in Europe, some of the many great mathematicians were the following:

- Bonaventura Cavalieri (1598–1647) in Italy,
- John Wallis (1616–1703) in England,
- Nicolaus Mercator (1620–1687) in Germany and England,
- Christian Huygens (1629–1695) in Holland,
- Isaac Barrow (1630–1677) in England,
- James Gregory (1638–1675) in Scotland,
- Isaac Newton (1642–1727) in England,
- Gottfried Wilhelm Leibniz (1646–1716) in Germany.

Cavalieri rediscovered Archimedes's *method of indivisibles* for calculating areas and volumes. As an example of Cavalieri's approach, let us consider the following derivation of the area of the ellipse $y = (b/a)(a^2 - x^2)^{\frac{1}{2}}$. First we plot the circle $y = (a^2 - x^2)^{\frac{1}{2}}$ on the same rectangular coordinate frame of reference. Second, we note that a vertical chord of the ellipse is just b/a of the corresponding vertical chord of the circle (i.e. the chord lying in the same straight line). Third, we think of the areas of the ellipse and the circle as somehow being the 'sums' of their chords. The ratio of ellipse chord to

corresponding circle chord is always b/a , so it follows that the ratio of the area of the ellipse to the area of the circle is also b/a . Hence the area of the ellipse is $(b/a)\pi a^2 = \pi ab$.

Cavalieri's most important result is perhaps the theorem which, in our notation, reads as follows:

$$\int_0^a x^n dx = \frac{a^{n+1}}{n+1}.$$

Wallis was a professor of geometry at Oxford. He was a Royalist, and ended up as chaplain to Charles II. He also invented a method for teaching deaf mutes.

In algebra, he used negative and fractional exponents: $x^{-n} = 1/x^n$ and $x^{p/q} = \sqrt[q]{x^p}$. Using Cavalieri's methods, he calculated the area under the curve

$$y = a_0x^0 + a_1x^1 + \cdots + a_nx^n.$$

He also discovered, but could not give a rigorous proof for, the following expression of π as an infinite product:

$$2 \cdot \frac{2}{1} \cdot \frac{2}{3} \cdot \frac{4}{3} \cdot \frac{4}{5} \cdot \frac{6}{5} \cdot \frac{6}{7} \cdot \frac{8}{7} \cdot \frac{8}{9} \cdots$$

Mercator calculated the area under the curve

$$y = \frac{1}{1+x} = 1 - x + x^2 - x^3 + \cdots,$$

obtaining

$$\log_e(1+x) = x - x^2/2 + x^3/3 - x^4/4 + \cdots$$

(converging for $-1 < x < 1$). (It was Gregory of St. Vincent (1584–1667) who first showed that the integral of $1/x$ is the natural log of x .) The cartographer's Mercator projection is due, not to Nicolaus Mercator, but to Gerhardus Mercator (1512–1592).

Like Pascal, Huygens made a study of the cycloid. This is the curve described by a point fixed to the rim of a wheel as it rolls along a flat surface. Huygens showed that the 'evolute' of a cycloid is again a cycloid. He also discovered (Galileo notwithstanding) that a pendulum had to swing in a cycloidal arc, and not in a circular arc, if the period of oscillation was to be strictly independent of the amplitude of the swing.

Huygen's most profound contribution to physics was the wave theory of light. Assuming that light is propagated in waves in an oscillating 'ether', he was able to explain the laws of reflection and refraction. His theory was in one respect better than Newton's corpuscular (particle) theory, which failed to explain interference phenomena. Nonetheless, Huygen's theory was

eclipsed by Newton's for many years. It is only in modern quantum mechanics that we acknowledge that light may be viewed as waves as well as particles.

Often a promising young mathematician is denied a job in mathematics because all the places are filled by older professors who have not done a thing since the day they were given tenure. These old professors might learn a lesson from Isaac Barrow. In 1669, he resigned his professorship of geometry at Cambridge, so that the young Newton might take his place. Thereafter, Barrow devoted himself to divinity.

Gregory used the method of Wallis to obtain 'Gregory's series':

$$x = \tan x - \frac{1}{3} \tan^3 x + \frac{1}{5} \tan^5 x - \cdots.$$

Putting $x = \pi/4$, this gives

$$\pi/4 = 1 - 1/3 + 1/5 - 1/7 + \cdots,$$

an interesting but slowly converging series for π . (This tan series had also been obtained by Indian mathematicians, and is found in a book called the *Tantrasangraha-vyakhya* (c. 1530).)

Newton never knew his father, who had died before his birth. Newton was supposed to learn farming, but spent his time doing experiments and building mechanical models. Finally, his uncle relented and sent him to Cambridge, where he read Euclid, Descartes, Kepler and Wallis, and attended the lectures of Isaac Barrow.

By the time he got his B.A., Newton had discovered derivatives, which he called 'fluxions', and had established the Binomial Theorem for integer and fractional exponents (without, however, giving a rigorous proof).

Newton used the Binomial Theorem to expand certain functions $f(x)$ into power series. When he wanted to find the area under the curve $y = f(x)$ he could then apply Wallis's method, replacing x^n by $x^{n+1}/(n+1)$. At this time, there was little knowledge about the conditions under which one can treat an infinite sum in the same way as a finite sum. Indeed, there is no evidence that Newton worried about the convergence of the series in his generalized Binomial Theorem.

During the plague of 1665–66, Newton withdrew to the family farm and thought about gravitation. Back in Cambridge, he first helped Barrow with some lecture notes, and then took over his professorship in 1669. Around this time he discovered the decomposition of white light by a prism. From 1673 to 1683 he lectured on algebra and the theory of equations.

Newton's greatest contribution to knowledge was his theory of universal gravitation, which once and for all provided a rational explanation for the

apparently erratic motions of the heavenly bodies. The fundamental idea was simple: the same force which causes an apple to fall must act on the moon and the planets. (Note that Newton had to reject the Aristotelian idea that an apple falls simply because it is the nature of 'earthy' things to go downwards.)

Similar ideas had previously occurred to Hooke, Huygens, Halley and Wren. These thinkers realized that Kepler's laws implied that any 'force of gravity' would have to obey an inverse square law, such as that given by Newton. What Newton did was to show that, conversely, an inverse square law implies Kepler's laws.

The *law of gravitation* asserts that, given two bodies with masses M and m , at a distance x apart, the force between them is given by

$$F = kMmx^{-2},$$

where k is a universal constant. If M is large compared with m , we may think of this force as being exerted by M on m . Since Newton had defined *force* to mean rate of change of momentum, in this case $F = -m\ddot{x}$, it follows that the acceleration $\ddot{x} = -kMx^{-2}$ of the smaller body does not depend on its mass m . This result is still valid today, even if Newton's law has to be slightly modified to conform to the general theory of relativity. If M is the mass of the earth, assumed to be concentrated at its center, and m is the mass of the apple on or near the surface, \ddot{x} is practically constant, confirming Galileo's original observation.

At first Newton worked out the planetary motions from the assumption that the sun and the planets were points. However, he was not happy about this assumption, and so he did not publish his results immediately. It was only in 1685, about twenty years later, that he was able to prove that the gravitational force due to a solid sphere is the same as if the entire mass were concentrated in the center. He assumed that the density of matter at a point inside the sphere depends only on its distance from the center, this presumably being the case with all the planets in our solar system.

Having overcome this last difficulty, Newton finally published his epoch-making *Principia* in 1686. To avoid all controversy about his methods, he replaced his original arguments involving the infinitesimal calculus by classical geometrical arguments in the style of Euclid, which his contemporaries were able to understand. Unfortunately, for this very reason, today the *Principia* is difficult to read.

The second volume of the *Principia* dealt with hydrostatics and hydrodynamics. It showed that Descartes's theory of vortices did not work. Newton's theories were soon accepted everywhere. Even in France, Voltaire advocated Newton against Descartes (in 1733).

It was only in 1692 that Newton published two letters on 'fluxions', as he called derivatives. He wrote \dot{x} , \ddot{x} for the first and second derivatives with respect to a parameter t (for time). He wrote o for dt and xo for dx .

In 1704, in an appendix to a book on optics, he gave a study of 'fluents', his name for indefinite integrals or antiderivatives. He showed the connection between fluents and 'quadratures', that is, definite integrals. He also treated maxima and minima, tangents to curves and lengths of curves.

In 1696, Newton abandoned his Cambridge professorship for a government position in London and, three years later, he was Master of the Mint.

We defer the discussion of Leibniz to the next chapter, ending this chapter with a quotation from Newton:

I do not know what I may appear to the world; but to myself I seem to have been only like a boy, playing on the sea-shore, and diverting myself, in now and then finding a smoother pebble, or a prettier shell than ordinary, whilst the great ocean of truth lay all undiscovered before me.

Exercises

1. Show that if the three Laws of Indices are to hold for negative and fractional exponents, Wallis's way of using such exponents is the only way. (The three Laws of Indices are $x^{m+n} = x^m x^n$, $(x^m)^n = x^{mn}$, $(xy)^m = x^m y^m$.)
2. Show that, when m is odd,

$$\int_0^{\pi/2} \sin^m x \, dx = \frac{(m-1)(m-3) \cdots 2}{m(m-2) \cdots 3 \cdot 1}.$$

Derive a similar formula for the case in which m is even. Finally, use the above two formulas to give a derivation of Wallis's product formula for π .

3. Suppose the wheel in the definition of the cycloid is rolling along the positive x axis, with the fixed point starting at the origin. What is the equation of the resulting cycloid if the radius of the wheel is 1?
4. How many terms of Gregory's series do you have to add up to get π accurate to two decimal places?
5. According to Newton's calculations, your weight on a planet of uniform density is determined only by the matter which is closer to the center than you are. Suppose that you dig down towards the center of such a planet (assumed spherical). Draw a graph showing how your weight varies as you descend towards the center. (Hint: at the center, your weight is 0.)

6. You are on a planet consisting of a spherical inner core of density 5 and radius R . This inner core is in the middle of a spherical outer core of radius $5R$ and density 1. Show that, as you dig down towards the center of this planet, your weight at first decreases, then increases, and, finally, decreases. Where does your weight start to increase?

29

Leibniz

Gottfried Wilhelm Leibniz (1646–1716) was born in Leipzig. His father died when he was six. Leibniz educated himself, using his late father's library, and entered the university in Leipzig when he was only fifteen. In 1666 he was refused the degree of Doctor of Law on the grounds that he was too young. In the same year, Leibniz conceived the idea of symbolic logic, a universal language in which all rational thinking could be expressed.

Leibniz worked as a diplomat for the Elector of Mainz. In this capacity, he went to Paris, where Louis XIV rejected his idea of attacking Egypt instead of another European country. Here Leibniz met Huygens, who introduced him to geometry and physics.

Huygens challenged Leibniz to sum the series

$$\sum_{n=1}^{\infty} \frac{2}{n(n+1)} = 1 + 1/3 + 1/6 + \cdots.$$

Leibniz solved the problem thus: $2/n(n+1) = 2(1/n - 1/(n+1))$, so the series equals $2(1 - 1/2 + 1/2 - 1/3 + 1/3 - 1/4 + 1/4 - \cdots) = 2(1 + 0) = 2$. In the 20th century, we would object to this on the grounds that Leibniz might equally well have written

$$2(2 - 3/2 + 3/2 - 4/3 + 4/3 - 5/4 + 5/4 - \cdots) = 2(2 + 0) = 4.$$

We now prefer to solve this problem by first showing that the m th partial sum of the series is $2(1 - 1/(m+1))$ and then taking the limit as $m \rightarrow \infty$.

In 1673, Leibniz visited England and became a fellow of the Royal Society. If he got his ideas about the Calculus from Newton, he never acknowl-

edged this. Two years later, he developed his own version of the Calculus, introducing the notation we still use today. In particular, he obtained the rule for the derivative of a product that is named after him.

Leibniz's work on the Calculus appeared in 1684, before Newton had got around to publishing his results. Not surprisingly, Leibniz was accused of plagiarism by some British mathematicians, not without Newton's acquiescence, and a bitter priority battle ensued. Today, we ascribe the invention of the Calculus to both Newton and Leibniz.

Leibniz had been working as librarian for the Duke of Hannover. When the latter became king of England as George I, he left two people behind: his wife, whom he divorced and shut up in a cloister, and Leibniz, because he did not want to antagonize the British academic establishment.

Leibniz thought of dy and dx in dy/dx as 'infinitesimals'. Thus dx was an infinitely small increment in x which was yet different from 0, and dy , defined as

$$dy = f(x + dx) - f(x)$$

for a given function $y = f(x)$, was also different from 0 (unless f happened to be constant near x). For example, if $f(x) = x^2$, then $dy = (x+dx)^2 - x^2 = 2x(dx) + (dx)^2$. This represented the 'rise' of the function corresponding to the 'run' dx . Hence, the slope of the tangent was rise/run = $dy/dx = 2x + dx$, so that, at x , the tangent has slope $2x$.

The concept of the infinitesimal — also implicit in Newton's fluxions — was criticized by many, including the philosopher and bishop George Berkeley (1685–1753). How, he asked, can we divide by dx if it is 0? How can we get the slope of the tangent right, in our example $2x$, if it is not 0?

Karl Weierstrass (1815–1897) agreed that there were problems, and he responded by putting calculus on the firm footing it has today. At the moment, dy/dx is not seen as a quotient but as a limit of quotients:

$$\frac{dy}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}.$$

If we insist on the separate existence of dx and dy , we may put $dx = h$ and $dy = \frac{dy}{dx}h$. Weierstrass also gave a perfectly rigorous definition of a limit (the *epsilon-delta* definition) which does not depend on vague notions like 'small', 'approaches', etc.

There is also at the moment a rigorous version of the infinitesimal itself. In 1966 Abraham Robinson introduced a *nonstandard model* for real numbers in which there is an entity ξ such that

$$0 < \xi < 1, \quad 0 < \xi < 1/2, \quad 0 < \xi < 1/3, \quad \dots \quad (*)$$

Can this be done consistently? Yes, and the reason, roughly speaking, is as follows. We know from mathematical logic that, if a contradiction were deducible from $(*)$, then this contradiction would be deducible in a finite

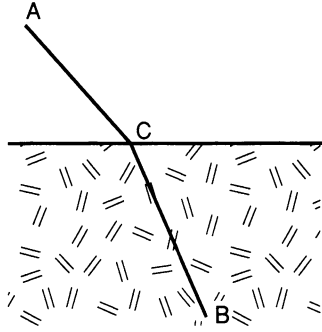


FIGURE 29.1. Refraction of light

number of steps, hence from a finite subset of (*). Since this is not possible (assuming that the ordinary interpretation of the rational numbers is consistent), there is no contradiction.

Leibniz's infinitesimals are simple and intuitive but not rigorous; Robinson's infinitesimals are rigorous but neither simple nor intuitive. Anyone studying nonstandard analysis will discover a complex and bizarre system. On the other hand, it should be stated that nonstandard analysis has helped some mathematicians obtain new results in ordinary analysis.

Leibniz also made contributions to philosophy. He is famous for his view that this is the best of all possible universes (Could God have failed to create the best?). This view was ridiculed by Voltaire in his novel *Candide*. However, it is unlikely that Leibniz intended 'best' to mean 'happiest'. What he had in mind may have been something like the discovery of Willebrord Snell, who was able to explain the refraction of light by assuming that a ray of light minimizes the time in going from point *A* in one medium to point *B* in another. (Think of a person trying to walk from point *A* on a paved surface to a point *B* in a rough field. Since he finds it harder to walk in the field, he will not go straight from *A* to *B*, but in a broken line *ACB*, as shown in Figure 29.1.) Leibniz makes explicit reference to Snell in Section XXII of the *Discourse on Metaphysics*, where Leibniz discusses the 'easiest way' in which a ray of light might travel. Thus it is quite possible that by 'best' universe, Leibniz meant, among other things, 'easiest', or 'most energy-efficient' universe.

The behaviour of light rays is only one instance of the *Principle of Least Action*, which was put forward by Pierre Maupertuis (1698–1759) in 1751. As now formulated, this principle asserts that any physical process happens in such a way as to minimize the action $\int_a^b E dt$, where *E* is the difference between kinetic and potential energy and *t* is the time. This principle was formally established by Lagrange.

Some utilitarian philosophers, such as J.S. Mill, have claimed that even human psychology is, or ought to be, determined by something like the

Principle of Least Action. In this they followed Plato, who had Socrates argue in *The Protagoras* that the way we ought to behave is to maximize pleasure minus pain, that virtue consists in knowing when to forego a present gain in favour of a larger future gain (temperance) and when to face the danger of an immediate loss against the expectation of ultimate gain (courage). However, Socrates seems to have had second thoughts about this later. When Crito urged him to flee from his prison so as to escape the death penalty and be able to bring up his sons, Socrates refused, believing it was his moral duty to remain in jail. He replied thus:

Whatever the popular view is, and whether the alternative is pleasanter than the present one or even harder to bear, the fact remains that to do wrong is in every sense bad and dishonourable (Plato's *Crito* 49b; see also *Laws* 707d).

Leibniz's most original contribution to philosophy is probably the system incorporated in his *Monadology*. In this work he proposes that the universe is made up of certain ultimate elements, called *monads*, which are capable of perception. Human souls are monads with memory and reason.

If Leibniz was influenced by his mathematical background in formulating this model of the universe, then he might have thought of a monad as a point together with the set of all points infinitely close to it. In fact, Robinson attributed this technical meaning to the word 'monad' in his nonstandard analysis.

On the other hand, Bertrand Russell thought of the monads as points, with arrows connecting different points (see his *Mysticism and Logic*). Each monad has a physical aspect, consisting of all arrows emerging from it, and a mental aspect, consisting of all arrows converging on it.

Leibniz asserted that a monad has no windows. (Its contact with the rest of the universe is via a 'pre-established harmony'.) This suggests that anything we know about a monad must be deducible from the 'arrows' relating it to all the other monads in the universe.

One is reminded here of the 20th century notion of a category (see Part II, Chapter 31). For instance, in the category of groups, we have as points all groups, and as arrows all group homomorphisms. If we want to know the elements of a group G , we need not look 'inside' the group at all, but only at the arrows from the free group in one generator to G .

Exercise

1. Suppose the half-plane $y > 0$ is a rough field where you would walk at u km/h, and the half-plane $y < 0$ is a paved surface where you would walk at v km/h. You are at point $(0, -6)$ and want to walk to the point (a, b) . How should you go to get there fastest?

The Eighteenth Century

We shall only mention some of the most important mathematicians of the eighteenth century:

- Brook Taylor (1685–1731),
- Colin Maclaurin (1689–1746),
- Abraham de Moivre (1667–1754),
- Leonhard Euler (1707–1783),
- Joseph Louis Lagrange (1736–1813),
- Pierre Simon Laplace (1749–1827),
- Adrien Marie Legendre (1752–1833).

Brook Taylor, an ardent admirer of Newton, discovered the *Taylor series*

$$f(a + x) = f(a) + xf'(a) + x^2 f''(a)/2! + \dots,$$

publishing it in 1715.

Colin Maclaurin, a Scotsman, is best known for the special case $a = 0$ of Taylor's series. This appeared in his *Treatise of Fluxions* (1742). In his book, Maclaurin tried to be sufficiently rigorous to answer Berkeley's objections to the Calculus, but he did not even get to the point of demonstrating conditions under which his *Maclaurin series* converges.

Abraham de Moivre was born in France, but lived in England. He is famous for his formula

$$(\cos x + i \sin x)^n = \cos nx + i \sin nx ,$$

which is easily proved, for natural numbers n , by mathematical induction. De Moivre published an important book on the theory of probability, called the *Doctrine of Chances*.

Society failed de Moivre: in spite of letters of recommendation from both Newton and Leibniz, he was never given a proper job in mathematics. He had to earn a meagre living by private tutoring and answering gamblers' questions on probability. It is said that, as he approached the end of his life, de Moivre slept fifteen minutes longer each day. When he reached a full twenty four hours, he died.

Although Leonhard Euler was Swiss, he spent part of his professional life in Berlin and most of it in St. Petersburg. Towards the end of his life he became blind, but this did not slow down his mathematical output. He found many interesting and exciting results in mathematics. Indeed, it has been said that Euler picked all the raisins out of the mathematical cake. Some of his results are the following:

1. If a convex polyhedron has V vertices, F faces and E edges, then $V + F - E = 2$. For example, a cube has 8 vertices, 6 faces and 12 edges; we have $8 + 6 - 12 = 2$. (Descartes came close to this formula, but he did not actually state it.)
2. $e^{i\pi} = -1$, where e is the 'Euler number':

$$e = \lim_{n \rightarrow \infty} (1 + 1/n)^n.$$

3. $1/1^2 + 1/2^2 + 1/3^2 + 1/4^2 + \cdots = \pi/6$.

Euler's proof of this was not rigorous but, before Euler, no one even guessed that the sum of the series was $\pi/6$.

4. Every even perfect number has the form $2^{n-1}(2^n - 1)$ where $2^n - 1$ is prime.
5. If n is a positive integer, let $\phi(n)$ be the number of natural numbers less than or equal to n and relatively prime to it. Then, if a is a positive integer relatively prime to n , it follows that n is a factor of $a^{\phi(n)} - 1$. Fermat's Little Theorem is a corollary of this.
6. The circumcenter, orthocenter and centroid of a triangle are collinear. The line that passes through them is called the *Euler line*.

(The *circumcenter*, *orthocenter* and *centroid* of a triangle are the meeting points of the right bisectors, altitudes and medians, respectively.)

7. Fermat was wrong when he conjectured that all natural numbers of the form $2^{2^n} + 1$ are primes.

Euler recognized the importance of convergence in dealing with infinite series, but he did not always pay attention to it. For example, he would write

$$1/(x-1) = 1/x + 1/x^2 + 1/x^3 + \cdots$$

(which is correct when $|x| > 1$) and put $x = 1/2$ to obtain $-2 = 2 + 4 + 8 + \cdots$. He also showed a lack of rigour in employing his principle of ‘conservation of form’, according to which a theorem true for natural number exponents also holds for any real exponent. In this way, he obtained facile ‘proofs’ of the generalized Binomial Theorem, and the generalized de Moivre’s Theorem.

In addition to his numerous discoveries in pure mathematics, by no means all of which have been discussed here, Euler also made important contributions to mechanics. He elaborated the Principle of Least Action. Finally, he worked out a theory of lunar motion. His collected works run to about 75 volumes.

The following proof of Euler’s formula $V + F - E = 2$ was suggested by H. S. M. Coxeter.

Let O be a point in the interior of the convex polyhedron. About O as center describe a sphere which contains the polyhedron. Now imagine a source of light placed at O . Rays emanating from O will project the polyhedron onto the surface of the sphere, mapping each flat polygon onto a spherical polygon whose sides are arcs of great circles. (This idea is said to be due to the Arabic mathematician Abu’l Wafa.) Choose a point in the interior of each spherical polygon and join it to the vertices by arcs of great circles, thus dividing each spherical polygon into as many spherical triangles as it has sides. Then

$$\begin{aligned} 720^\circ &= \text{area of sphere} \\ &= \text{sum of areas of spherical triangles} \\ &= \text{sum of angles of spherical triangles} - 180^\circ \times \text{number of triangles} \\ &= \text{sum of angles at interior points} + \\ &\quad \text{sum of angles at vertices} - 180^\circ \times 2E \\ &= 360^\circ \times F + 360^\circ \times V - 360^\circ \times E. \end{aligned}$$

Dividing by 360° , we obtain Euler’s formula.

Joseph Lagrange was born in Italy of mixed French and Italian parentage. His father lost the family fortune through speculation, but Lagrange later commented that, if it had not been for this, he might never have turned to mathematics. He was converted to mathematics through an essay by Halley.

At age 23, Lagrange was able to explain, on the basis of Newton's theory of gravitation, why the moon always shows the same face to the earth.

Having acquired an early fame, Lagrange spent 25 years in Prussia at the invitation of Frederick II. After Frederick's death, Lagrange moved to Paris, where he became a favourite of Marie Antoinette. He had mixed feelings about the Revolution, especially when his friend, the chemist Lavoisier, was guillotined, but he stuck it out. He was involved in the introduction of the decimal system for weights and measures. When people pleaded the advantages of the base 12, he would ironically defend the base 11. He became professor of mathematics at the Ecole Polytechnique.

Lagrange was a universal mathematical genius, his interests ranging from number theory to physics. Among his achievements are the following:

1. The first proof of Wilson's Theorem that, if p is a prime number, then it is a factor of $(p-1)! + 1$;
2. The first complete solution of the Diophantine equation $x^2 - Ry^2 = 1$, where R is a given nonsquare positive integer; Lagrange generalized this to give a complete treatment of Diophantine equations of the form

$$Ax^2 + Bxy + Cy^2 + Dx + Ey + F = 0,$$

where A, B, C, D, E , and F are given integers;

3. The first proof that every natural number is a sum of four squares of natural numbers (e.g., $7 = 2^2 + 1^2 + 1^2 + 1^2$ and $9 = 3^2 + 0^2 + 0^2 + 0^2$);
4. A systematic theory of differential equations;
5. The *Mécanique*, which he conceived at the age of 19 but only published at 52, in which he expressed the dynamics of a rigid system by the equations

$$\frac{d}{dt} \frac{\partial T}{\partial \dot{\theta}} - \frac{\partial T}{\partial \theta} + \frac{\partial V}{\partial \theta} = 0,$$

where T is the total kinetic energy, V is the potential energy, t is the time, θ is any coordinate, and $\dot{\theta} = d\theta/dt$; Lagrange observed that his equations expressed the fact that the total action $\int_a^b (T - V)dt$ was minimal; to justify this observation, he had to invent the calculus of variations.

When Lagrange wrote to Euler about his results in the calculus of variations, Euler was so impressed that he withheld his own results from publication so that Lagrange could publish first. Sad to say, such unselfish acts are rare.

After his wife died in 1783, Lagrange wore himself out publishing the *Mécanique*. The excesses of the Revolution upset him and he became subject to fits of depression. From these the lonely genius was rescued by the love of a teenaged girl, Renée Le Monnier, who insisted on marrying him in 1792. For the remaining twenty years of his life, Lagrange was both happy and mathematically productive.

Laplace was the son of humble parents but ended up as a marquis under the restored Bourbons. Politically, he was an opportunist, but occasionally he stood up for his principles. Napoleon once told him, 'you have written a big book on the universe without mentioning its creator', to which Laplace replied: 'I don't need that hypothesis'.

Laplace was more of a mathematical physicist than a pure mathematician. He introduced the potential V and showed that it satisfied

$$\frac{\partial^2 V}{\partial x^2} + \frac{\partial^2 V}{\partial y^2} + \frac{\partial^2 V}{\partial z^2} = 0.$$

His greatest contribution to mathematics was the useful phrase 'it is easy to see', which peppers his *Mécanique Céleste*. In *The History of Mathematics*, David Burton reports:

The American astronomer Nathaniel Bowditch (1773-1838), who translated four of the five volumes into English, observed, "I never came across one of Laplace's 'Thus it plainly appears' without feeling sure that I had hours of hard work before me to fill up the chasm and find out and show how it plainly appears."

Legendre was a great promoter of Euclid. He showed that the Parallel Postulate follows from the assumption that the plane contains real squares (i.e., quadrilaterals with four equal sides, each of whose angles is a right angle.) He also did work on the method of least squares.

Legendre is best known for his work in number theory. He was the first to prove that the Diophantine equation

$$x^5 + y^5 = z^5$$

has no nonzero integer solutions. He introduced the Legendre symbol $\left(\frac{n}{p}\right)$, where p is a prime and n an integer not divisible by p . He wrote $\left(\frac{n}{p}\right) = 1$ when n has the form $kp + r^2$ (with k and r integers) and $\left(\frac{n}{p}\right) = -1$ when n does not have this form.

Euler had conjectured a theorem, called the Law of Quadratic Reciprocity. Using the Legendre symbol, it can be expressed as follows: if p and q are distinct odd primes, $\left(\frac{p}{q}\right)\left(\frac{q}{p}\right) = (-1)^{(p-1)(q-1)/4}$. Euler could not prove this, but Legendre found proofs for some special cases. It was Gauss who published the first complete proof, in 1801. Gauss later gave five other proofs of the same result.

The Legendre symbol $\left(\frac{n}{p}\right)$ can be calculated quite easily, in view of the following observation due to Euler:

$$n^{\frac{p-1}{2}} = \left(\frac{n}{p}\right) + \text{a multiple of } p.$$

Indeed, it follows from Fermat's Little Theorem that

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

is a multiple of p , so p divides either $n^{\frac{p-1}{2}} - 1$ or $n^{\frac{p-1}{2}} + 1$ (but not both, since it does not divide their difference). We claim that p divides the former if and only if $\left(\frac{n}{p}\right) = 1$, hence p divides the latter if and only if $\left(\frac{n}{p}\right) = -1$, from which facts the observation follows.

To see this, at least in one direction, suppose $\left(\frac{n}{p}\right) = 1$, that is, $n = r^2 + kp$ for some integers r and k . Then

$$n^{\frac{p-1}{2}} = (r^2 + kp)^{\frac{p-1}{2}} = r^{p-1} + k'p = 1 + k''p$$

for some integers k' and k'' , hence p divides $n^{\frac{p-1}{2}} - 1$. The converse implication, though not difficult, is a little tricky, and we shall omit its proof.

Exercises

1. Give conditions sufficient for the convergence of the Maclaurin series.
2. Prove de Moivre's formula for positive integers n .
3. Show that e , as defined above, is bound below by 2 and above by 3.
4. Prove that the circumcenter, orthocenter and centroid of any triangle are collinear.
5. Give an example of a formula with exponents which is true when the exponents are natural numbers but not always true when the exponents are rationals.
6. Prove Wilson's Theorem.
7. Check the Law of Quadratic Reciprocity for $p = 5$ and $q = 13$.

The Law of Quadratic Reciprocity

In this chapter we single out a great 19th century mathematician and present one of his most elegant proofs. The reader should be warned, however, that the 19th century was so rich in mathematics that it really deserves a book of its own.

Carl Friedrich Gauss (1777–1855) was inspired to become a mathematician by his discovery of a ruler and compass construction for the regular polygon with 17 sides — this when he was only a teenager — but his gift had revealed itself much earlier: as a three year old, he had pointed out an error in his father's payroll accounts!

Gauss's first major contribution was his Number Theory book *Disquisitiones Arithmeticae*, which appeared in 1801. As well as the first construction for the regular 17-gon, it included the first proof of the Fundamental Theorem of Arithmetic (that every integer can be uniquely written as a product of primes), the first proof that every prime p has a primitive root g (meaning that no two of the numbers $1, g, g^2, \dots, g^{p-2}$ differ by a multiple of p), the first proof that every natural number is a sum of three triangular numbers, and the first proof of the theorem featured below, namely, the Law of Quadratic Reciprocity.

Gauss was also an astronomer and a physicist. In 1807, subsequent to his calculation of the position of the asteroid Ceres, he was appointed director of the observatory at Göttingen and, in 1809, he published a book on planetary astronomy. In physics, Gauss did pioneering work in electromagnetism; the *gauss* is a unit of measure denoting magnetic intensity.

We conclude this chapter with a short proof of the quadratic reciprocity law, adapted from one of the proofs given by Gauss.

Let p and q be distinct odd primes and write $f(x, y) = py - qx$. Consider the region of the Cartesian plane consisting of all pairs of real numbers (x, y) such that

$$1/2 < x < p/2, \quad 1/2 < y < q/2.$$

This region is subdivided into four mutually disjoint parts

$$A : f(x, y) < -q/2, \quad B : -q/2 < f(x, y) < 0,$$

$$C : 0 < f(x, y) < p/2, \quad D : p/2 < f(x, y).$$

Replacing x by $\frac{p+1}{2} - x$ and y by $\frac{q+1}{2} - y$, we establish a one-to-one correspondence between A and D . For any subset S of the Cartesian plane, let $L(S)$ denote the number of pairs of integers (x, y) in S ; then clearly $L(A) = L(D)$. Now

$$L(A) + L(B) + L(C) + L(D) = \frac{p-1}{2} \cdot \frac{q-1}{2},$$

so that

$$L(B) + L(C) = \frac{p-1}{2} \cdot \frac{q-1}{2} - 2L(A).$$

Thus, the law of quadratic reciprocity can be written

$$\left(\frac{p}{q}\right) \left(\frac{q}{p}\right) = (-1)^{L(B)+L(C)},$$

and this will follow if we show that

$$\left(\frac{p}{q}\right) = (-1)^{L(B)}, \quad \left(\frac{q}{p}\right) = (-1)^{L(C)}.$$

For reasons of symmetry, it will suffice to prove one of these equations, say the latter. In view of how the Legendre symbol is calculated, this may be stated in the following form, known as Gauss's Lemma:

$$q^{\frac{p-1}{2}} = (-1)^{L(C)} + \text{a multiple of } p.$$

To prove Gauss's Lemma, consider all pairs of integers (x, y) in the following region:

$$E : 1/2 < x < p/2, \quad 1/2 < y < q/2, \quad -p/2 < f(x, y) < p/2.$$

In E ,

$$y - 1/2 < \frac{qx}{p} < y + 1/2,$$

so that y is the integer *closest to* $\frac{qx}{p}$; we shall write $y = g(x)$ to indicate that y is completely determined by x . Clearly, E contains C ; in fact, $(x, y) \in C$ if and only if $(x, y) \in E$ and $f(x, y) > 0$.

As x ranges from 1 to $\frac{p-1}{2}$, and therefore y from 1 to $\frac{q-1}{2}$, the absolute value of $f(x, y) = f(x, g(x))$ takes exactly $\frac{p-1}{2}$ values. For, if $|f(x, y)| = |f(x', y')|$, where x' and y' range as x and y , then

$$q(x \pm x') = p(y \pm y'),$$

hence $x \pm x' = kp$ for some integer k , so that $x = x'$ and therefore $y = y'$.

Since, in E , $|f(x, y)| < p/2$, these $\frac{p-1}{2}$ values are $1, \dots, \frac{p-1}{2}$ and hence

$$\prod_{x=1}^{\frac{p-1}{2}} |f(x, y)| = \left(\frac{p-1}{2}\right)!$$

Finally, we calculate

$$q^{\frac{p-1}{2}} \left(\frac{p-1}{2}\right)! = \prod_{x=1}^{\frac{p-1}{2}} qx = \prod_{x=1}^{\frac{p-1}{2}} (py - f(x, g(x))).$$

This differs by a multiple of p from

$$\prod_{x=1}^{\frac{p-1}{2}} (-f(x, g(x))) = (-1)^{L(C)} \prod_{x=1}^{\frac{p-1}{2}} |f(x, g(x))| = (-1)^{L(C)} \left(\frac{p-1}{2}\right)!,$$

since $-f(x, g(x))$ will be negative if and only if $(x, g(x)) \in C$. Since p does not divide $\left(\frac{p-1}{2}\right)!$, it follows that $q^{\frac{p-1}{2}}$ differs from $(-1)^{L(C)}$ by a multiple of p . This completes the proof of Gauss's Lemma and therefore the proof of the Law of Quadratic Reciprocity.

Part II

Topics in the Foundations of Mathematics

1

The Number System

Much of 19th and 20th century mathematics is not accessible or meaningful at the undergraduate level. Still, we plan to examine some important and up to date material, including some exciting recent discoveries.

The following letters are used to represent sets of numbers:

N the natural numbers $0, 1, 2, \dots$;

Z the integers (Z being for ‘Zahlen’);

Q the rationals (Q being for ‘quotients’);

R the reals;

C the complex numbers $a + bi$;

H the quaternions (H being for Hamilton, the person who introduced them).

This now widely accepted notation was first proposed by N. Bourbaki, actually a slowly changing group of French mathematicians who have been engaged for half a century in writing up the *Elements of Mathematics* in a systematic and trend setting fashion. One of the founding members was the late Jean Dieudonné, to whom many of their early decisions may be attributed.

The number systems are arranged here in what mathematicians conceive to be the correct logical order, which differs from the historical one (see the Introduction). Zero was not originally considered to be a natural number

and positive rationals and reals were studied before negative integers were considered.

Each of the sets after \mathbf{N} 'extends' the preceding one and, with the exception of \mathbf{H} , is motivated by our desire to solve equations which are otherwise unsolvable. Thus, \mathbf{Z} , \mathbf{Q} , \mathbf{R} , and \mathbf{C} are needed to solve the equations $x + 2 = 0$, $2x = 3$, $x^2 = 2$ and $x^2 + 1 = 0$, respectively.

Mathematicians often 'construct' a set on this list from the set just above it. For example, they build the rationals out of the integers by equating rationals with certain classes of pairs of integers: the fraction $2/3$ is equated with $\{(2, 3), (4, 6), \dots, (-2, -3), \dots\}$. However, the list begins with the natural numbers. How were they constructed? According to Kronecker, 'God created the natural numbers, all the other numbers were made by man'. Still, as we shall see, some people have tried to construct the natural numbers also. This has led one wit to suggest that 'man created the natural numbers, all the others were Dieudonné'.

We should note that there are other kinds of numbers that do not appear on our list, such as the transfinite numbers and the infinitesimals.

What are the natural numbers? In particular, what is the number 2? We know it is not just the sign or numeral '2' (or 'II'). It is not something perishable or changing. But what is it? There is more than one answer.

The *Platonists* hold that numbers are abstract, necessarily existing objects. The number two is that Platonic 'form' or 'idea' in virtue of which things have the property of two-ness.

For the *logicists*, numbers are things which can be defined in terms of logic. For example, for Bertrand Russell, 2 is

$$\{x \mid \exists_y \exists_z (y \neq z \wedge x = \{y, z\})\}.$$

In other words, 2 is the set of all unordered pairs. On the other hand, von Neumann claimed that 2 is the set $\{0, 1\}$, where $1 = \{0\}$, and 0 is just the empty set. Pursuing yet another approach, Church held that 2 is the process of iteration

$$\lambda_f(f \circ f),$$

more precisely, the mapping which assigns to every function f its iterate $f \circ f$, defined by $(f \circ f)(x) = f(f(x))$.

For the *formalists*, 2 is just a class of expressions manipulated according to certain rules. Often they do not define numbers, but rather give axioms that characterize them. For example, Peano sees 2 as $SS0$, where ' $SS0$ ' is a string of symbols which are manipulated according to certain axioms (see below).

The *intuitionists* hold that numbers are mental entities which would not exist unless people thought about them. For Brouwer, 2 is the concept which expresses the principle of 'two-ity'.

We see that each of these schools has a different view of the matter. What is the true answer? Fortunately, Professor Tournesol was able to discover this when attending a conference in France. The number 2 is a pair of platinum balls kept at room temperature in the second drawer at the Bureau of Standards in Paris.

2

Natural Numbers (Peano's Approach)

For Peano, the natural number system is a triple $(N, 0, S)$, where N is a set, 0 is an element of that set, and S is a function whose domain and codomain are that set, such that the following axioms hold for all elements x and y of \mathbf{N} :

1. $Sx \neq 0$,
2. $Sx = Sy \Rightarrow x = y$,
3. when $\phi(x)$ is any proposition involving the natural number x , if
 - (a) $\phi(0)$,
 - (b) $\forall_{x \in N} \phi(x) \Rightarrow \phi(Sx)$,

then $\phi(x)$ is true for any natural number x .

We define 1 as $S0$ and 2 as $SS0 = S1$, etc. (Normally, we leave out the parentheses.)

Note that (3) above is not a single axiom, but a whole ‘scheme’ or class of axioms, one for each ϕ . This scheme is called *mathematical induction*.

In Peano’s system we define addition as follows:

$$\begin{aligned}x + 0 &= x, \\x + Sy &= S(x + y).\end{aligned}$$

Such definitions are called *recursive* definitions. We add 3 and 2 thus:
 $3 + 2 = 3 + S1 = S(3 + 1) = S(3 + S0) = S(S(3 + 0)) = SS(SSS0) = 5$.

Multiplication is defined recursively as follows:

$$\begin{aligned}x \cdot 0 &= 0, \\x \cdot Sy &= (x \cdot y) + x.\end{aligned}$$

We sometimes write $(x \cdot y) + x$ as $xy + x$. Assuming that we already know how to add, we can now multiply. Thus $3 \cdot 2 = 3 \cdot S1 = (3 \cdot 1) + 3 = (3 \cdot S0) + 3 = ((3 \cdot 0) + 3) + 3 = (0 + 3) + 3 = 3 + 3 = 6$.

Exponentiation is defined recursively as follows:

$$\begin{aligned}x^0 &= 1, \\x^{Sy} &= (x^y) \cdot x.\end{aligned}$$

Here even $0^0 = 1$. The reader who has already learned how to multiply will now have no difficulty in calculating 3^2 . All the usual laws of arithmetic follow from the above axioms and definitions; for example:

commutativity: $x + y = y + x$, $xy = yx$;

associativity: $x + (y + z) = (x + y) + z$, $x(yz) = (xy)z$;

distributivity: $x(y + z) = xy + xz$, $(x + y)z = xz + yz$;

laws of indices: $(x^y)^z = x^{zy}$, $x^y x^z = x^{y+z}$, $x^z y^z = (xy)^z$.

We give some examples of how the above laws (and others) follow from Peano's axioms and definitions.

Lemma 2.1. $0 + x = x$.

Proof: We prove this using mathematical induction. First we show that it is true when $x = 0$: $0 + 0 = 0$ by the definition of addition (df+). Second we assume it holds for x and prove it for Sx . Our assumption that it is true for x is called the *induction hypothesis* (hyp). We have

$$\begin{aligned}0 + Sx &= S(0 + x) && \text{(df+)} \\ &= S(x + 0) && \text{(hyp)} \\ &= Sx && \text{(df+)}.\end{aligned}$$

Since the assumption that the result holds for x implies that it holds for Sx , and since it is true for 0, we may conclude, by mathematical induction, that it is true for all natural numbers.

Lemma 2.2. $Sx + y = S(x + y)$.

Proof: Use induction on y . The result is true when $y = 0$: $Sx + 0 = Sx$ (df +) and $Sx = S(x + 0)$ (df +). If the result holds for y ,

$$\begin{aligned}Sx + Sy &= S(Sx + y) && \text{(df+)} \\ &= SS(x + y) && \text{(hyp)} \\ &= S(x + Sy) && \text{(df+)}.\end{aligned}$$

Theorem 2.3. $x + y = y + x$.

Proof: By induction on y . The result follows from Lemma 2.1 when $y = 0$. Supposing the result holds for y , we have

$$\begin{aligned} x + Sy &= S(x + y) && (\text{df } +) \\ &= S(y + x) && (\text{hyp}) \\ &= Sy + x && (\text{Lemma 2.2}). \end{aligned}$$

Theorem 2.4. $x + (y + z) = (x + y) + z$.

Proof: Use induction on z :

$$x + (y + 0) = x + y = (x + y) + 0.$$

Supposing the result holds for z , we argue thus:

$$\begin{aligned} x + (y + Sz) &= x + S(y + z) && (\text{df } +) \\ &= S(x + (y + z)) && (\text{df } +) \\ &= S((x + y) + z) && (\text{hyp}) \\ &= (x + y) + Sz && (\text{df } +). \end{aligned}$$

Theorem 2.5. $x^z y^z = (xy)^z$.

Proof: $x^0 y^0 = 1 \cdot 1 = 1 \cdot S0 = (1 \cdot 0) + 1 = 0 + 1 = 1 = (xy)^0$, so the result is true when $z = 0$. Suppose it holds for z . Then $x^{Sz} y^{Sz} = (x^z x)(y^z y)$ and $(xy)^{Sz} = (xy)^z xy = (x^z y^z)xy$ (hyp). The result now follows by mathematical induction on z — provided we can first establish the associativity and commutativity of multiplication. This we leave to the reader.

Theorem 2.6. If $x + y = x + z$ then $y = z$.

Proof: By Theorem 2.3, this is true when $x = 0$. Assume it holds for x . If $Sx + y = Sx + z$ then $S(x + y) = S(x + z)$ (Lemma 2.2) and hence, by Peano's second axiom, $x + y = x + z$. By the induction hypothesis, $y = z$.

Every natural number except 0 has a predecessor. We can define a *naive predecessor* as follows:

$$\begin{aligned} P0 &= 0, \\ PSy &= y. \end{aligned}$$

Given the naive predecessor function, it is easy to define naive subtraction. Again, we use a recursive definition:

$$\begin{aligned} x \dot{-} 0 &= x, \\ x \dot{-} Sy &= P(x \dot{-} y). \end{aligned}$$

The reason we use the sign $\dot{-}$ rather than the sign $-$ is that naive subtraction is not quite the same as ordinary subtraction. We cannot say that $1 - 3 = -2$ since -2 is not a natural number. Instead we say that $1 \dot{-} 3 = 0$. Using the above definition, we have

$$\begin{aligned} 1 \dot{-} 3 &= 1 \dot{-} S2 = P(1 \dot{-} 2) = P(1 \dot{-} S1) = P(P(1 \dot{-} 1)) = P(P(1 \dot{-} S0)) \\ &= P(P(P(1 \dot{-} 0))) = P(P(P(1))) = P(P(P(S0))) = P(P(0)) = P(0) = 0. \end{aligned}$$

We define $\min(x, y)$ as $x \dot{-} (x \dot{-} y)$, and $\max(x, y)$ as $x + y \dot{-} \min(x, y)$.

Giuseppe Peano (1858–1932) published essentially this system in his *Arithmetices Principia* (1889), a book written in a language he invented.

Exercises

1. Prove that $0 \cdot x = 0$ without using the commutative law for multiplication.
2. Prove $1^x = 1$.
3. Prove $x(y + z) = xy + xz$.
4. Prove $xy = yx$.
5. Prove $x \dot{-} x = 0$.
6. Prove $Sx \dot{-} Sy = x \dot{-} y$.
7. Prove $\max(x, y) + \min(x, y) = x + y$.
8. In the *Arithmetices Principia*, Peano actually defines $x \dot{-} y$ as

$$x \dot{-} y = \begin{cases} z & \text{such that } y + z = x \\ 0 & \text{if there is such a } z \\ & \text{otherwise} \end{cases}$$

- (a) Why can't there be two natural numbers z and w such that $y + z = x$ and $y + w = x$?
- (b) Show that Peano's definition of $\dot{-}$ is equivalent to the recursive definition given above.

3

The Integers

It is not difficult, though rather boring, to construct the integers from the natural numbers. Instead, we shall demonstrate in the next chapter how to construct the rationals from the integers, by essentially the same process. But first let us state the properties which make the set of integers into what is called an *integral domain*.

A *ring* $(R, 0, -, +, 1, \cdot)$ is a set R with operations $0, -, +, 1$, and \cdot , where $+$ and \cdot are *binary* operations, $-$ is a *unary* operation and 0 and 1 are *nullary* operations, that is, specified elements of R , which moreover satisfy the following axioms or identities:

1. $(x + y) + z = x + (y + z)$ (associativity),
2. $x + 0 = x$,
3. $x + (-x) = 0$,
4. $x + y = y + x$ (commutativity),
5. $x \cdot 1 = x = 1 \cdot x$,
6. $(x \cdot y) \cdot z = x \cdot (y \cdot z)$ (associativity),
7. $(x + y) \cdot z = (x \cdot z) + (y \cdot z)$,
 $z \cdot (x + y) = (z \cdot x) + (z \cdot y)$ (distributivity).

We may take this opportunity to review a bit of abstract algebra. Axioms

(1) to (3) make $(R, 0, -, +)$ into a group. It is not difficult to prove that, in a group,

$$(2') \quad 0 + x = 0,$$

$$(3') \quad -x + x = 0.$$

A group is said to be *Abelian* if it also satisfies the commutative law (4). We have stated this as an axiom, even though it is a consequence of the remaining axioms of a ring (not of a group). The operation \cdot may or may not obey the commutative law for multiplication:

$$(8) \quad x \cdot y = y \cdot x.$$

If it does, we call the ring *commutative*. For example, \mathbf{Z} and \mathbf{Q} are commutative rings, but the ring of 2×2 matrices with entries from \mathbf{Z} is not commutative.

A commutative ring is called an *integral domain* if

$$(9) \quad 0 \neq 1 \text{ and } x \cdot y = 0 \text{ implies } x = 0 \text{ or } y = 0.$$

It is called a *field* provided

$$(10) \quad 0 \neq 1 \text{ and, if } x \neq 0, \text{ there exists an element } y \text{ such that } x \cdot y = 1 = y \cdot x.$$

\mathbf{Z} is an integral domain, but not a field. On the other hand, \mathbf{Q} is a field as well as an integral domain. In fact, every field is an integral domain.

Exercises

1. Prove that (2') and (3') must hold in a group.
2. Prove that, in a ring, $x \cdot 0 = 0 = 0 \cdot x$ and $x \cdot (-1) = -x = (-1) \cdot x$.
3. Prove that (4) follows from the other axioms of a ring.

4. Show that every field is an integral domain.
5. In a group, if $0'$ is another element such that, for all x , $x + 0' = x$, show that $0' = 0$.
6. In a group, if \sim is another unary operation such that, for all x , $x + (\sim x) = 0$, show that $\sim x = -x$.
7. Prove that, in any integral domain, we have the following *cancellation law*: if $a \cdot c = b \cdot c$ and $c \neq 0$ then $a = b$.

4

The Rationals

In this chapter we shall construct the field of rationals from the ring of integers. Our argument will use only properties (1) to (9) of the integers (Part II, Chapter 3), and hence the *same* argument will produce a field from *any* integral domain.

One's first idea is to say that a rational number a/b is a pair of integers (a, b) , where $b \neq 0$. But this does not work, because then $2/3 \neq 4/6$. So we shall begin by introducing an equivalence relation on pairs of integers (a, b) such that $b \neq 0$. We define

$$(a, b) \equiv (c, d) \iff ad = bc,$$

it being assumed, of course, that b and d are nonzero. This is an *equivalence relation*, that is

$$(a, b) \equiv (a, b) \quad \text{(reflexivity),}$$

$$\text{if } (a, b) \equiv (c, d), \text{ then } (c, d) \equiv (a, b) \quad \text{(symmetry),}$$

$$\text{if } (a, b) \equiv (c, d) \text{ and } (c, d) \equiv (e, f), \text{ then } (a, b) \equiv (e, f) \quad \text{(transitivity).}$$

These three properties of an equivalence relation are easily checked; we shall verify only transitivity.

Given $ad = bc$ and $cf = de$, we wish to show that $af = be$. Making free

use of associativity and commutativity, we calculate

$$afd = adf = bcf = bde = bed.$$

As we are given that $d \neq 0$, we may use the cancellation law (Exercise 7 of Chapter 3) to infer that $af = be$, as was required to be shown.

Assuming that $b \neq 0$, we define the *ratio* (rational number) a/b as the equivalence class of (a, b) , that is, as the set of all pairs (c, d) , with $d \neq 0$, such that $ad = bc$. Note that

$$a/b = c/d \iff (a, b) \equiv (c, d) \iff ad = bc.$$

If \mathbf{Q} is the set of such ratios, we shall obtain a field $(\mathbf{Q}, \underline{0}, -, +, \underline{1}, \cdot)$ where we have underlined $\underline{0}$ and $\underline{1}$ to distinguish them, at least temporarily, from the integers 0 and 1.

Here is how we attempt to define the operations of \mathbf{Q} :

$$\begin{aligned}\underline{0} &= 0/1, \\ -(a/b) &= (-a)/b, \\ a/b + c/d &= (ad + bc)/(bd), \\ \underline{1} &= 1/1, \\ a/b \cdot c/d &= (ac)/(bd).\end{aligned}$$

Note that, because of axiom (9) for an integral domain, $bd \neq 0$. There is no problem with $\underline{0}$ and $\underline{1}$, but we must check that the other operations are *well-defined*. We shall do so here for the operation $+$.

Thus, suppose that $a/b = a'/b'$ and $c/d = c'/d'$. We must show that

$$a/b + c/d = a'/b' + c'/d',$$

that is,

$$(ad + bc)/(bd) = (a'd' + b'c')/(b'd'),$$

that is,

$$(ad + bc)(b'd') = (bd)(a'd' + b'c').$$

The reader is invited to verify that this is indeed the case.

The ratios with operations $\underline{0}$, $-$, $+$, $\underline{1}$, and \cdot as defined above form a field. To prove this one must check all the axioms of a field. For example, we shall check (4) and (10). To prove (4) for ratios we argue as follows:

$$\begin{aligned}a/b + c/d &= (ad + bc)/(bd) \\ &= (da + cb)/(db) \text{ (commutativity of } \cdot \text{ for integers)} \\ &= (cb + da)/(db) \text{ (commutativity of } + \text{ for integers)} \\ &= c/d + a/b.\end{aligned}$$

Here we made use of the commutativity of \cdot and $+$ for integers. To prove (10), we may assume that $a/b \neq \underline{0}$, that is, $a1 \neq b0$, that is, $a \neq 0$ (see Exercise 2 of Chapter 3). We claim that $(a/b) \cdot (b/a) = 1$. Indeed,

$$(a/b) \cdot (b/a) = (a \cdot b)/(b \cdot a) = 1/1,$$

since

$$(a \cdot b) \cdot 1 = a \cdot b = b \cdot a = (b \cdot a) \cdot 1.$$

We have thus constructed the field \mathbf{Q} from the ring \mathbf{Z} . More generally, the same construction leads from any integral domain to a field, called its *field of quotients*.

What is the relationship between \mathbf{Q} and the ring \mathbf{Z} from which it is constructed? Strictly speaking, the set \mathbf{Q} does not contain the set \mathbf{Z} . However, \mathbf{Q} does contain a subset, consisting of the ratios of the form $a/1$, which is *isomorphic* to \mathbf{Z} . To make this notion more precise, consider the mapping $h : \mathbf{Z} \rightarrow \mathbf{Q}$ such that $h(a) = a/1$. Then h is a *homomorphism*, that is, it preserves the operations of \mathbf{Z} :

$$\begin{aligned} h(a+b) &= h(a) + h(b), \\ h(ab) &= h(a)h(b), \\ h(-a) &= -h(a), \\ h(0) &= \underline{0}, \\ h(1) &= \underline{1}. \end{aligned}$$

Furthermore, h is a 1-to-1 mapping (an *injection*). The set $h(\mathbf{Z})$ is thus an *isomorphic image* of \mathbf{Z} . We also say that h *embeds* \mathbf{Z} in \mathbf{Q} . It is a mathematical convention to identify 2 and $2/1$ and to say that \mathbf{Z} is contained in \mathbf{Q} .

Exercises

1. Show that $-(-(a/b)) = a/b$.
2. Show that $((a/b)^{-1})^{-1} = a/b$ (assuming $a, b \neq 0$).
3. Prove that $a/(b/c) = (ac)/b$.
4. Define $\max(a/b, c/d)$.
5. What definition might one use to construct the integers out of the natural numbers?
6. If b is a nonzero integer, let $b\mathbf{Z} = \{bx | x \in \mathbf{Z}\}$. Consider any mapping $f : b\mathbf{Z} \rightarrow \mathbf{Z}$ such that $f(bx + by) = f(bx) + f(by)$ for all $x, y \in \mathbf{Z}$. Show that $f(0) = 0$ and $f(-bx) = -f(bx)$. If g is any other function

like f (that is, $g : c\mathbf{Z} \rightarrow \mathbf{Z}$ for $c \neq 0$ and $g(cx + cy) = g(cx) + g(cy)$) define $f \equiv g$ to mean that $f(bcx) = g(bcx)$ for all integers x . Show that \equiv is an equivalence relation and that the equivalence class $[f]$ can be taken as a definition of the quotient $f(b)/b$.

5

The Real Numbers

There are two well-known ways of constructing the reals from the rationals: the Dedekind cut approach, which goes back to Eudoxus, and the Cauchy sequence approach.

Eudoxus, a member of Plato's Academy, was interested in defining the notion of proportion for geometric quantities. What he did can be interpreted in modern terms as defining *equality* between two real numbers (ratios of geometric quantities) α and β as follows: $\alpha = \beta$ iff the set of all rationals below α is the same as the set of all rationals below β , and similarly for the sets of rationals above α and β . Dedekind exploited this idea further by defining the real number α as the pair (L, U) of sets of rationals below and above α , respectively. Such pairs can be described without mentioning α ; e.g., L might be the set of all rationals x for which $x^2 < 2$ and U the set of rationals y for which $y^2 > 2$. We shall not develop this idea further, as it is discussed in many algebra books.

Analysts prefer a construction of the reals proposed by Cauchy, according to which the real number α is defined as the set of all sequences of rational numbers which converge to α . Again, this can be done without mentioning α . A *Cauchy sequence* is a sequence $\{a_n | n \in \mathbf{N}\}$ of rational numbers such that $|a_m - a_n|$ can be made as small as one likes by taking m and n sufficiently large. Two Cauchy sequences $\{a_n | n \in \mathbf{N}\}$ and $\{b_n | n \in \mathbf{N}\}$ are said to be *equivalent* if $|a_n - b_n|$ can be made as small as one likes by taking n sufficiently large. A *real number* is then defined as an equivalence class of Cauchy sequences.

An amusing construction of the real numbers, which is not so well-known, bypasses the rationals altogether and defines a real number as an equiva-

lence class of certain mappings $f : \mathbf{Z} \rightarrow \mathbf{Z}$. Call f *almost linear* if the set of all $|f(m+n) - f(m) - f(n)|$, where $m, n \in \mathbf{Z}$, is bounded. Call two almost linear mappings $f, g : \mathbf{Z} \rightarrow \mathbf{Z}$ *equivalent* if the set of all $|f(n) - g(n)|$, where $n \in \mathbf{Z}$, is bounded. Define a *real number* to be an equivalence class of almost linear mappings $\mathbf{Z} \rightarrow \mathbf{Z}$. The real number α will be the equivalence class of the mapping f for which $f(n) = [\alpha n]$ is the greatest integer $\leq \alpha n$. Addition and multiplication of the real numbers corresponding to the almost linear mappings f and g are easily defined as the equivalence classes of $f + g$ and $f \circ g$, respectively, where $(f + g)(n) = f(n) + g(n)$ and $(f \circ g)(n) = f(g(n))$. This definition is due to Steve Schanuel.

Instead of constructing the real numbers, one may describe the field of real numbers axiomatically as a *complete ordered field*. We already know what is meant by a field, so we only have to define the words ‘ordered’ and ‘complete’.

A field F is *ordered* if it has a subset P such that

1. $x, y \in P \Rightarrow x + y \in P$,
2. $x, y \in P \Rightarrow xy \in P$,
3. exactly one of the following holds: $x = 0$, or $x \in P$, or $-x \in P$.

Note that by (3) either 1 or -1 is an element of P . Since $(-1)(-1) = 1$, it follows from (2) that $1 \in P$. The elements of P are the *positive* elements of the field. The existence of P allows us to define an *order relation* on the field F :

$$x \leq y \iff x = y \text{ or } y - x \in P.$$

From this definition we obtain the following propositions:

$x \leq x$	reflexivity,
$x \leq y$ and $y \leq z \Rightarrow x \leq z$	transitivity,
$x \leq y$ and $y \leq x \Rightarrow x = y$	antisymmetry,
$x \leq y$ or $y \leq x$	dichotomy.

It is because \leq has these four properties that it is called an *order relation*. \mathbf{Q} and \mathbf{R} are ordered fields, but \mathbf{C} is not. To see that \mathbf{C} is not ordered, consider the element i . If $i \in P$ then $-1 = i \cdot i \in P$. This is impossible since $1 \in P$. If $-i \in P$ we again have that $-1 = (-i)(-i) \in P$. So neither i nor $-i$ is in P , and this contradicts (3) above.

An ordered field is *complete* if every nonempty set of positive elements has a greatest lower bound. For example, $\sqrt{2}$ is the greatest lower bound of all positive reals r such that $2 \leq r^2$. Since $\sqrt{2}$ is not rational, we may conclude from this example that \mathbf{Q} is not a complete ordered field. There is really only one complete ordered field in the sense of the following:

Theorem 5.1. *Any two complete ordered fields are isomorphic.*

For the proof, see Chapter 4.3 of Birkhoff and Mac Lane [1977].

Exercises

1. Pick any of the three definitions of the reals mentioned here and prove that they form a field.
2. Prove the four properties of the order relation defined above.

6

Complex Numbers

Negative numbers are required to solve the equation $x + 2 = 1$. To solve $2x = 3$, we need the rationals. The equation $x^2 = 2$ has an irrational solution. Finally, the equation $x^2 = -1$ has an *imaginary root* called i .

Numbers of the form $a + bi$ where a and b are real, are called *complex numbers*. Complex numbers with $a = 0$ are also called *imaginary*. The complex number $a + bi$ is often associated with the point (a, b) in the Cartesian plane. The absolute value of a complex number is just its distance $\sqrt{a^2 + b^2}$ from the origin. The angle θ measured counterclockwise from the positive x axis to the line joining (a, b) to the origin is called the *angle* of the complex number $a + bi$; thus $\tan \theta = b/a$.

Complex numbers were introduced by Girolamo Cardano (1501–1576), who used them in his *Ars Magna* (1545) to solve cubic equations. Cardano tells us to multiply $5 + \sqrt{-15}$ by $5 - \sqrt{-15}$, ‘putting aside the mental tortures involved’ (see T. Richard Witmer’s translation of the *Ars Magna*, p. 219).

There are many ways to define the complex numbers, all of them being essentially equivalent. For example, we can define them as ordered pairs of real numbers subject to the multiplication rule

$$(a, b) \cdot (c, d) = (ac - bd, ad + bc).$$

Here $(1, 0)$ plays the role of 1 and $(0, 1)$ the role of i .

Another way to introduce the ‘field’ of complex numbers is to say that it is the quotient ring $\mathbf{R}[x]/(x^2 + 1)$. The elements of this ring are equivalence classes of polynomials with real coefficients, where two polynomials are said

to be equivalent if they differ by a multiple of $x^2 + 1$. Here the role of i is played by the equivalence class whose representative is the polynomial x .

Yet another way to define complex numbers is to say that they are 2×2 matrices with real entries of the form

$$\begin{pmatrix} a & b \\ -b & a \end{pmatrix}.$$

It is easy to see that the set of matrices of this form is closed under addition and multiplication. Here the role of 1 is played by the identity matrix and the role of i is played by

$$\begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}.$$

The advantage of defining complex numbers in this way is that one can use the arithmetic of matrices to give a quick proof of the fact that the complex numbers form a ring. The commutative law of multiplication still needs checking, since it does not hold for matrices in general. Moreover, the inverse of a nonzero matrix of the form

$$\begin{pmatrix} a & b \\ -b & a \end{pmatrix}$$

is just the matrix

$$\begin{pmatrix} a/k & -b/k \\ b/k & a/k \end{pmatrix},$$

where $k = a^2 + b^2$ is the determinant of the former.

Since a real number is not usually thought of as a matrix, one might well ask how the reals relate to the complex numbers if the latter are conceived as matrices. Well, the mapping $h: \mathbf{R} \rightarrow \mathbf{C}$, where

$$h(a) = \begin{pmatrix} a & 0 \\ 0 & a \end{pmatrix}$$

is a 1-to-1 homomorphism. Hence $h(\mathbf{R})$ is an isomorphic image of \mathbf{R} and may be identified with \mathbf{R} , and we may say that \mathbf{R} is a subset of \mathbf{C} .

Let u and v be any complex numbers. Then

$$|u + v| \leq |u| + |v|,$$

$$|uv| = |u| |v|.$$

If we think of the complex numbers as points in the Cartesian plane, the above inequality is just the triangle inequality of Euclid (Book I, 20). The above equality is equivalent to the algebraic identity

$$(ac - bd)^2 + (ad + bc)^2 = (a^2 + b^2)(c^2 + d^2).$$

This identity was known to al-Khazini in 950 AD. (See Part I, Chapter 23.)

Associating $a + bi$ with the point (a, b) in the Cartesian plane, let r be its distance from the origin, and let θ be its angle. Then

$$a + bi = r(\cos \theta + i \sin \theta).$$

If $a' + b'i$ has absolute value r' and angle θ' , we may exploit the well-known addition formulas for trigonometric functions to obtain

$$(a + bi)(a' + b'i) = rr'(\cos(\theta + \theta') + i \sin(\theta + \theta')).$$

By induction on the natural number m , this formula leads to the equation

$$(r(\cos \theta + i \sin \theta))^m = r^m(\cos(m\theta) + i \sin(m\theta)),$$

which is known as de Moivre's theorem after Abraham de Moivre (1667–1754), who was the first to make use of it (Part I, Chapter 30).

Exercises

1. Using the matrix definition of complex numbers, verify that \mathbf{C} is a field.
2. What is the multiplicative inverse of $5 + \sqrt{-15}$?
3. Give an algebraic proof of the triangle inequality.
4. Prove the theorem of de Moivre.

The Fundamental Theorem of Algebra

The negative integers, fractions, irrationals and imaginary numbers were all introduced in order to supply solutions to polynomial equations. Do we need more numbers? Is there a polynomial equation with complex coefficients with a root that is not a complex number? The answer is no, and the fact that no new numbers are required is called the Fundamental Theorem of Algebra. It was first stated by Albert Girard (1595–1632) in 1629, and proofs were given by Jean d’Alembert (1717–1783) and Carl Friedrich Gauss (1777–1855). However, John Stillwell [1989] argues that there is a flaw in Gauss’s proof, and the first rigorous proof was given only after Weierstrass established the basic properties of continuous functions.

Theorem 7.1. (The Fundamental Theorem of Algebra)

Every polynomial equation with complex coefficients has complex solutions.

Proof. We assume some simple facts about continuity. For example, we assume that if a continuous closed curve which loops around the origin n times gradually changes so that it loops around the origin only $n - 1$ times, then, at some time, part of it passes through the origin. Let

$$w = p(z) = z^m + c_1 z^{m-1} + \cdots + c_{m-1} z + c_m \quad (c_m \neq 0),$$

where the c_k ’s are complex numbers. Without loss of generality, we may take our polynomial equation to be $p(z) = 0$. Let

$$g(z) = \frac{p(z)}{z^m} = 1 + c_1/z + \cdots + c_m/z^m.$$

Then

$$\begin{aligned}
 |g(z) - 1| &= |c_1/z + \cdots + c_m/z^m| \\
 &\leq |c_1/z| + \cdots + |c_m/z^m| \\
 &= \frac{|c_1|}{|z|} + \cdots + \frac{|c_m|}{|z|^m} \\
 &\leq m \left(\frac{\max(|c_k|)}{|z|} \right)
 \end{aligned}$$

if $|z| > 1$. Geometrically, this means that, when $|z|$ is large, $g(z)$ is represented by a point close to $(1, 0)$. In particular, if z moves around a large circle with centre at the origin, then $g(z)$ will move in some continuous closed curve near $(1, 0)$ and not loop around $(0, 0)$.

We consider two planes, one to represent $z = x + yi$, and one to represent $w = u + vi$. As z moves around a circle of radius $|z| = r$ and centre $(0, 0)$, w moves in some continuous closed curve in its own plane.

Suppose that z moves around a circle with radius large enough to keep $g(z)$ from winding around $(0, 0)$. How many times does w loop around $(0, 0)$? Equivalently, how does the angle of w change as z moves around its large circle? In the previous chapter we noted that the angle of a product of two complex numbers is the sum of the angles of those two numbers. Since $w = z^m g(z)$, the angle of w changes by an amount equal to the change in the angle of z^m plus the change in the angle of $g(z)$. Since $g(z)$ stays close to $(1, 0)$ the net change in the angle of $g(z)$ as z goes around the large circle is 0. Thus the change in the angle of w is just the change in the angle of z^m . But the angle of a product of complex numbers is equal to the sum of their angles, so the change in the angle of w is just m times the change in the angle of z . As z moves around a circle, the change in the angle of z is 2π . Hence the change in the angle of w is $2\pi m$. In other words, as z moves around the large circle, w loops around the origin m times.

Now imagine the radius of the large circle gradually decreasing to 0. When it is close to 0, w moves in a continuous closed curve close to $p(0) \neq 0$. (We are assuming that $c_m \neq 0$.) Thus w no longer loops around the origin; at some point as the radius of the large circle decreased, w moved through the origin, that is, $p(z)$ was equal to 0 and the polynomial equation had a complex number solution.

Corollary 7.2. *A polynomial of degree m with complex coefficients has exactly m linear factors.*

Proof: We have seen that $p(z) = 0$ has a complex number solution, call it z_1 . If $p(z) = q(z)(z - z_1) + R$ then $R = p(z_1) = 0$, so $z - z_1$ is a factor of $p(z)$. Now $q(z)$ has degree $m - 1$, and it also has a complex root. Continuing the process, we can write $p(z) = (z - z_1)(z - z_2) \cdots (z - z_m)$.

It is a pity that the proof of the fundamental theorem of algebra makes use of the notion of continuity, which belongs to analysis rather than to algebra. It often happens in mathematics that, when we want to prove a basic property of a certain system, we have to go outside that system. Still, the above proof can be employed constructively for solving polynomial equations. Before the advent of modern computers, there was a mechanical device which traced out the w curve for any given radius $|z| = r$. As r was gradually reduced, a bell would ring as soon as the w curve passed through the origin and a solution had been found.

Exercises

1. Let $w = z^2 - z - 2$. Graph w as z moves around the origin on a circle of radius 2.1.
2. Repeat the above exercise with a circle of radius 2.
3. Assuming that the coefficients of $p(z)$ are all real numbers, show that if $x + yi$ is a root of $p(z)$ then so is $x - yi$.

8

Quaternions

It was William Rowan Hamilton (1805–1865) who first conceived complex numbers as ordered pairs of reals, subject to the multiplication rule $(a, b)(c, d) = (ac - bd, ad + bc)$. His greatest contribution to pure mathematics, however, was his creation of an algebraic system in which the commutative law of multiplication does not hold. Just as many people before Bolyai and Lobachevsky had thought of Euclid's fifth postulate as a necessary and sacred truth, so many people before Hamilton believed that the law of commutativity for the multiplication of numbers was decreed by heaven. Hamilton discovered that there are consistent algebraic systems for which this law does not hold. (Matrix algebra came fourteen years later.)

What Hamilton discovered were 'quaternions'. The idea dawned on him while he was strolling along the Royal Canal in Dublin in 1843.

Quaternions are numbers of the form

$$a + bi + cj + dk,$$

where a, b, c , and d are real, and $i^2 = j^2 = k^2 = ijk = -1$.

If quaternion multiplication were commutative, it would be easy to derive $ij = -ij$ and from this it would follow that $i = j = 0$, and the whole system would collapse. We must have noncommutative multiplication. In particular, $ij \neq ji$. How do we know that such entities as i, j , and k exist? There are in fact three complex matrices:

$$i_1 = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}, \quad i_2 = \begin{pmatrix} 0 & i \\ i & 0 \end{pmatrix}, \quad i_3 = \begin{pmatrix} i & 0 \\ 0 & -i \end{pmatrix}$$

such that $i_1^2 = i_2^2 = i_3^2 = i_1 i_2 i_3 = -1$, where 1 is interpreted as the identity matrix. We may define the quaternion $a = a_0 1 + a_1 i_1 + a_2 i_2 + a_3 i_3$ as the complex matrix

$$\begin{pmatrix} a_0 + ia_3 & a_1 + ia_2 \\ -a_1 + ia_2 & a_0 - ia_3 \end{pmatrix},$$

namely, $\begin{pmatrix} u & v \\ -\bar{v} & \bar{u} \end{pmatrix}$, where $u = a_0 + ia_3$ and $v = a_1 + ia_2$, the a_k being real numbers. Addition and multiplication of quaternions is just the usual matrix arithmetic. The *conjugate* of a quaternion, as distinguished from the conjugate of the complex matrix which represents it, is the quaternion

$$\bar{a} = a_0 1 - a_1 i_1 - a_2 i_2 - a_3 i_3.$$

Note that we often identify the real number a_0 with the matrix $a_0 1$. Using this identification, we have

$$a\bar{a} = a_0^2 + a_1^2 + a_2^2 + a_3^2 = \bar{a}a.$$

The *norm* $N(a)$ of a quaternion a is the product $a\bar{a}$ viewed as a real number, namely, $a_0^2 + a_1^2 + a_2^2 + a_3^2$. Note that the norm is 0 just in case the quaternion is the zero matrix. For a nonzero quaternion a , we can write

$$a(\bar{a}/N(a)) = 1 = (\bar{a}/N(a))a.$$

(Here $N(a)$ is the real number, not the matrix.) This means that nonzero quaternions have multiplicative inverses.

A *division ring* is a system of elements closed under two binary operations, addition and multiplication, such that

1. under addition, it is a commutative group,
2. under multiplication, the nonzero elements form a group, and
3. multiplication distributes over addition (on both sides).

A division ring is thus a ring (Chapter 3) in which every nonzero element has an inverse under multiplication. It satisfies all the axioms of a field, except the commutative law of multiplication. It is sometimes called a *skew-field*. A key fact about quaternions is that they form a division ring.

It is also possible to represent quaternions by 4×4 real matrices. A cheap way to obtain such a representation is to replace i by $\begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}$ and 1 by $\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ in the complex 2×2 matrix representing the quaternion. Instead we shall look at a more interesting way to obtain such a representation, which has wider implications.

The quaternion $x = x_0 + x_1i_1 + x_2i_2 + x_3i_3$ gives rise to a column vector

$$[x] = \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix}.$$

Given a quaternion a , consider the mapping which assigns to every such vector $[x]$ the vector $[ax]$. Clearly, this is a linear transformation of the four-dimensional vector space over \mathbf{R} . Hence there will be a 4×4 matrix $A = L(a)$ such that $A[x] = [ax]$. We say that $L(a)$ *represents* the quaternion a . In particular, $I = L(1)$ is the identity matrix. For example, the matrix $I_1 = L(i_1)$ is obtained as follows: the equation

$$i_1(x_0 + x_1i_1 + x_2i_2 + x_3i_3) = -x_1 + x_0i_1 - x_3i_2 + x_2i_3$$

may be interpreted thus:

$$\begin{pmatrix} 0 & -1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} -x_1 \\ x_0 \\ -x_3 \\ x_2 \end{pmatrix}.$$

Hence $I_1 = L(i_1)$ is the 4×4 matrix appearing above. The matrices $I_2 = L(i_2)$ and $I_3 = L(i_3)$ are obtained similarly. Note that

$$L(ab)[x] = [(ab)x] = [a(bx)] = L(a)[bx] = L(a)L(b)[x]$$

for all vectors $[x]$. Therefore

$$L(ab) = L(a)L(b),$$

from which it easily follows that L is an injective homomorphism from the division ring of quaternions into the ring of all 4×4 matrices over \mathbf{R} .

Exercises

1. Show that the conjugate of ab is $\bar{b}\bar{a}$.
2. Show that $N(\bar{a}) = N(a)$ and $N(ab) = N(a)N(b)$.
3. Prove that every quaternion satisfies a quadratic equation with real coefficients.
4. Find all quaternions whose square is 1.
5. Prove that $L(\bar{a}) = A^t$ is the transpose of $A = L(a)$ and that $(N(a))^2 = \det(A)$, the determinant of A .

6. Show that with every quaternion a one may associate another 4×4 real matrix $R(a)$ such that $R(a)[x] = [xa]$ for all quaternions x .
7. Prove that $R(ab) = R(b)R(a)$ and $L(a)R(b) = R(b)L(a)$ for all quaternions a and b .

Quaternions Applied to Number Theory

In this chapter we shall use integer quaternions to show that every natural number is a sum of four perfect squares. This was first proved by J. L. Lagrange in 1770 (*Oeuvres*, Vol. 3, pp. 189-201).

As a warming up exercise, note that every integer is a sum of five cubes. Indeed, let m be an integer. Since $m - m^3 = -(m-1)m(m+1)$, it follows that $m - m^3$ is divisible by both 2 and 3, and hence by 6. Thus $x = (m - m^3)/6$ is an integer. Moreover, $m = m^3 + 6x = m^3 + (x+1)^3 + (x-1)^3 + (-x)^3 + (-x)^3$, a sum of five cubes. It is not known whether every integer can be written as a sum of four cubes.

About sums of non-negative cubes, it is known that every natural number, except 23 and 239, can be written as a sum of 8 non-negative cubes. As of 1971, it was not known whether the 8 could be lowered for large positive integers (Ellison [1971], pp. 10-36).

To prove the theorem of Lagrange, we shall require the following lemma, due to Euler.

Lemma 9.1. *For every odd prime p there exist integers x and y such that*

$$x^2 + y^2 + 1 = mp,$$

where m is an integer such that $0 < m < p$.

Proof: Let x range from 0 to $\frac{1}{2}(p-1)$. The squares x^2 all leave different remainders when divided by p . For suppose x_1^2 and x_2^2 leave the same remainder. Then $(x_1 + x_2)(x_1 - x_2) = x_1^2 - x_2^2$ is a multiple of p , hence p must divide $x_1 + x_2$ or $x_1 - x_2$. Without loss of generality, we may assume that

$x_1 > x_2$. Then $x_1 \neq x_2$ and

$$0 < x_1 + x_2 < p - 1, \quad -\frac{p-1}{2} \leq x_1 - x_2 \leq \frac{p-1}{2},$$

hence p divides neither $x_1 + x_2$ nor $x_1 - x_2$. Thus we have a contradiction and the assertion has been proved.

Similarly, we can show that, as y ranges from 0 to $\frac{1}{2}(p-1)$, the numbers $-y^2 - 1$ all leave different remainders when divided by p .

As x and y range from 0 to $\frac{1}{2}(p-1)$, the set of all x^2 thus takes on $\frac{1}{2}(p+1)$ different values and so does the set of all $-y^2 - 1$. Since there are only p possible remainders when one divides by p , the two sets must overlap; hence there exist integers x and y in the given range such that $x^2 + y^2 + 1 = mp$ is a multiple of p . Moreover,

$$1 \leq mp \leq \frac{1}{4}(p-1)^2 + \frac{1}{4}(p-1)^2 + 1 < p^2,$$

hence $1 \leq m < p$, as required.

Following Lipschitz [1886], p. 404, we define an *integer quaternion* as a quaternion with integer coefficients.

Theorem 9.2. (Lagrange)

Every natural number n is the sum of four perfect squares, that is, n is the norm of an integer quaternion.

Proof: Since the norm of the product of integer quaternions is the product of their norms and since n is a product of primes, it suffices to show that every prime is the norm of an integer quaternion. Since $2 = 1^2 + 1^2 + 0^2 + 0^2$, it suffices to prove this for odd primes. Let p be any odd prime. Then we know from Euler's lemma that there is an integer quaternion x such that $N(x) = mp$ with $0 < m < p$. Pick $m = m_0$ as small as possible with this property. We claim that $m_0 = 1$.

First let us show that m_0 cannot be even. If it is, then so is $x_0^2 + x_1^2 + x_2^2 + x_3^2$, hence also $x_0 + x_1 + x_2 + x_3$ is even. There are three cases: either all the x_i are even, or they are all odd, or exactly two are even, say x_0 and x_1 . In all three cases, $x_0 \pm x_1$ and $x_2 \pm x_3$ are even, hence

$$\frac{1}{2}m_0 = \left(\frac{x_0 + x_1}{2}\right)^2 + \left(\frac{x_0 - x_1}{2}\right)^2 + \left(\frac{x_2 + x_3}{2}\right)^2 + \left(\frac{x_2 - x_3}{2}\right)^2$$

is the sum of four perfect squares. But $\frac{1}{2}m_0$ is a positive integer less than m_0 , which contradicts the assumption that m_0 was chosen as small as possible.

We now know that m_0 is odd. Let z_i be the closest integer to $\frac{x_i}{m_0}$, hence $|\frac{x_i}{m_0} - z_i| < \frac{1}{2}$. (It cannot be equal to $\frac{1}{2}$, or else $m_0 = 2|x_i - m_0 z_i|$ would be even.)

Consider the integer quaternion $y = x - m_0 z$, where $z = z_0 + z_1 i_1 + z_2 i_2 + z_3 i_3$. Then

$$|y_i| = |x_i - m_0 z_i| < \frac{1}{2} m_0,$$

hence $N(y) < 4(\frac{1}{2} m_0)^2 = m_0^2$. But

$$N(y) = y\bar{y} = x\bar{x} - m_0(x\bar{z} + \bar{z}x) + m_0^2 z\bar{z}.$$

Write $x\bar{z} = w$, so $x\bar{z} + \bar{z}x = 2w_0$, where w_0 is the scalar part of w , and hence

$$N(y) = m_0 p - 2m_0 w_0 + m_0^2 N(z) = m_0 m_1,$$

where $m_1 = p - 2w_0 + m_0 N(z)$. Now $m_0 m_1 = N(y) < m_0^2$, hence $m_1 < m_0$.

Consider now the integer quaternion

$$y\bar{x} = x\bar{x} - m_0 z\bar{x} = m_0 p - m_0 z\bar{x} = m_0(p - z\bar{x}).$$

Then

$$m_0 m_1 m_0 p = N(y)N(\bar{x}) = N(y\bar{x}) = m_0^2 N(p - z\bar{x}),$$

hence

$$m_1 p = N(p - z\bar{x}).$$

Since $m_1 < m_0$, this would contradict the assumption that m_0 was chosen as small as possible, unless $m_1 = 0$.

This leaves only the possibility that $m_1 = 0$, hence $N(y) = 0$, hence $y = 0$, hence $x = m_0 z$, hence $m_0 p = N(x) = m_0^2 N(z)$, hence $p = m_0 N(z)$, hence $m_0 = 1$ or $m_0 = p$. But $m_0 < p$, so $m_0 = 1$, as was required.

Exercises

1. Express 239 as a sum of nine positive cubes.
2. Show that numbers of the form $8k + 7$ cannot be expressed as sums of three perfect squares.
3. Prove that every prime number of the form $4k + 1$ can be expressed as the sum of two perfect squares. (Hint: imitate the above proof using complex integers instead of integer quaternions.)

Quaternions Applied to Physics

If Hamilton's intention had been to apply quaternions to physics, he was stymied by the fact that physical space has only three dimensions. Writing the quaternion x as $x = x_0 + \xi$, where $\xi = i_1x_1 + i_2x_2 + i_3x_3$, we call x_0 the *scalar part* and ξ the *vector part*. A 3-vector is then a quaternion whose scalar part is 0. Unfortunately, if we multiply ξ by another vector $\eta = y_1i_1 + y_2i_2 + y_3i_3$, we obtain not a vector, but the quaternion

$$-(x_1y_1 + x_2y_2 + x_3y_3) + (y_2x_3 - x_2y_3)i_1 + \cdots$$

Oliver Heaviside pointed out the importance of the two separate parts of this product and wrote

$$\xi \circ \eta = x_1y_1 + x_2y_2 + x_3y_3,$$

$$\xi \times \eta = (y_2z_3 - x_2y_3)i_1 + \cdots$$

calling them the *scalar* and *vector product*, respectively. His *vector analysis* soon replaced the use of quaternions in physics. It enabled him to give a concise formulation of Maxwell's laws of electromagnetism. Writing $\nabla = i_1 \frac{\partial}{\partial x_1} + i_2 \frac{\partial}{\partial x_2} + i_3 \frac{\partial}{\partial x_3}$ for the vector representing partial differentiation with respect to the space coordinates and letting $E = E_1i_1 + E_2i_2 + E_3i_3$ and $M = M_1i_1 + M_2i_2 + M_3i_3$ represent the electric and magnetic fields, respectively, Maxwell's equations may be written as follows:

$$\nabla \circ M = 0, \quad \nabla \times E + \frac{\partial M}{\partial t} = 0, \quad \nabla \circ E = \rho, \quad \nabla \times M - \frac{\partial E}{\partial t} = \rho \frac{d\xi}{cdt},$$

where c is the velocity of light, ρ is the charge density, and $\frac{d\xi}{dt}$ the velocity of the matter bearing the electric charge.

Using the language of quaternions, albeit quaternions with complex components, these four equations may be combined into one, namely

$$\left(\frac{\partial}{c\partial t} - i\nabla\right)(M + iE) + \left(\rho + i\rho\frac{d\xi}{cdt}\right) = 0.$$

This was pointed out by Silberstein [1924] pp. 44-46; but one may wonder whether it wasn't already known to Maxwell himself, in view of his assertion that 'the invention of the calculus of quaternions is a step towards the knowledge of quantities related to space which can only be compared, for its importance, with the invention of triple coordinates by Descartes. The ideas of this calculus... are fitted to be of the greatest use in all parts of science.' (See Maxwell [1869].)

Einstein's theory of relativity made it clear that space and time should be combined into a single entity and that the expression

$$s^2 = c^2t^2 - x^2 - y^2 - z^2$$

should be invariant under coordinate transformations passing from a stationary to a moving platform. The minus sign in this expression suggests that we are talking about the norm of a quaternion

$$x = x_0 + i\xi, \quad x_0 = ct, \quad \xi = i_1x_1 + i_2x_2 + i_3x_3,$$

whose scalar part is real, but whose vector part is imaginary. Such a quaternion represents a point in so-called *Minkowski space*. Following Silberstein [1924], pp. 154-55, we observe that $s^2 = N(x)$ is invariant under a *Lorentz transformation*, which may itself be expressed with the help of 'biquaternions'.

A *biquaternion* $a = a_0 + i_1a_1 + i_2a_2 + i_3a_3$ has complex components a_0, \dots, a_3 . We must here distinguish the quaternion conjugate a^t of a from its complex conjugate a^c :

$$a^t = a_0 - i_1a_1 - i_2a_2 - i_3a_3,$$

$$a^c = \bar{a}_0 + i_1\bar{a}_1 + i_2\bar{a}_2 + i_3\bar{a}_3.$$

If a is represented as a 4×4 matrix $L(a)$ over \mathbf{C} as in Chapter 8, $L(a^t) = L(a)^t$, the transpose of $L(a)$. We call a biquaternion x *Hermitian* if $x^t = x^c$; this is what characterizes the quaternion $x = ct + i\xi$ describing a point in Minkowski space, that is, an event in space-time. A *Lorentz transformation* sends x onto pxp^{ct} , where p is a biquaternion of norm 1. Indeed,

$$(pxp^{ct})^t = p^c x^t p^t = p^c x^c p^t = (pxp^{ct})^c,$$

so pxp^{ct} is also Hermitian. Moreover, $N(pxp^{ct}) = pxp^{ct}p^c x^t p^t = pxx^t p^t = N(x)N(p) = N(x)$, since $p^{ct}p^c = (p^t p)^c = N(p)^c = 1^c = 1$. Thus a Lorentz transformation preserves the norm of a Hermitian biquaternion.

Another Hermitian biquaternion which transforms like x is

$$m_0 \frac{dx}{ds} = m + im \frac{d\xi}{cdt},$$

where $m = m_0 \frac{cdt}{ds}$ is the mass of the moving particle and $m \frac{d\xi}{dt}$ is its momentum. Here m_0 is called the *rest mass*; it is assumed to be invariant under Lorentz transformations. As Einstein observed, the principle of conservation of momentum should carry with it also that of

$$m = m_0 \frac{cdt}{ds} = m_0(1 - v^2/c^2)^{-\frac{1}{2}},$$

since $(\frac{ds}{dt})^2 = c^2 - v^2$, where $v^2 = (\frac{dx_1}{dt})^2 + (\frac{dx_2}{dt})^2 + (\frac{dx_3}{dt})^2$ is the square of the velocity of the particle. If v is small compared to c , we have

$$mc^2 = m_0 c^2 (1 - v^2/c^2)^{-\frac{1}{2}} \approx m_0 c^2 + \frac{1}{2} m_0 v^2,$$

which must then also be conserved. Einstein considered this to be the total energy $E = mc^2$ of the particle, consisting of the rest energy $m_0 c^2$ and the kinetic energy $\frac{1}{2} m_0 v^2$.

As we have seen, Maxwell's equations may be condensed into

$$\left(\frac{\partial}{c\partial t} - i\nabla \right) (M + iE) + \left(\rho + i\rho \frac{d\xi}{cdt} \right) = 0.$$

This may be written more concisely:

$$D^c F + J = 0,$$

where $D = \frac{\partial}{c\partial t} + i\nabla$ is sent by a Lorentz transformation onto $p^c D p^t$ and the so-called *six-vector* $F = M + iE$ is sent onto $p^c F p^{ct}$, so that $J = \rho + i\rho \frac{d\xi}{cdt}$ is sent onto $p J p^{ct}$. Thus J transforms like the mass-momentum biquaternion and may also be written as $\rho_0 \frac{dx}{ds}$, where

$$\rho = \rho_0 \frac{cdt}{ds} = \rho_0 \left(1 - \frac{v^2}{c^2} \right)^{-\frac{1}{2}}.$$

If the above rules for transforming D, F and J are adopted, it thus becomes clear that Maxwell's equations are preserved under Lorentz transformations. This fact, though without the aid of quaternions, appears to have first been noted by Poincaré.

Quaternions in Quantum Mechanics

What about quantum mechanics? If we adopt the representation of quaternions by 2×2 complex matrices (Chapter 8), the matrices ii_1 , ii_2 and ii_3 are known as *Pauli spin matrices*, except for sign.

Now let us consider the relativistic form of Schrödinger's wave equation for the electron. This is known as the *Klein-Gordon equation* and is usually written

$$\left(\frac{\partial^2}{c^2 \partial t^2} - \nabla \circ \nabla \right) \phi = -\mu^2 \phi,$$

where $\mu = 2\pi m_0/h$ is proportional to the rest-mass m_0 of the electron, h being Planck's constant. Using biquaternion notation, we write this

$$D^c D\phi = -\mu^2 \phi.$$

It is assumed that $\phi = \phi_0 + i\phi_1$ is a complex valued function of the position x in Minkowski space.

The second order Klein-Gordon equation may be replaced by two first order equations as follows. Putting $D\phi = \mu\chi$, where $\chi = \chi_0 + i\chi_1$, we obtain

$$\mu D^c \chi = D^c D\phi = -\mu^2 \phi.$$

Hence the Klein-Gordon equation is equivalent to the following pair of equations:

$$D\phi = \mu\chi, \quad D^c \chi = -\mu\phi.$$

Can these be combined into one first order equation?

Assume for the moment that there is an entity j such that $j^2 = -1$, $ji = -ij$, and $ji_k = i_kj$ for $k = 1, 2$ and 3 . Then we have

$$\begin{aligned} D(\phi + j\chi) &= D\phi + Dj\chi \\ &= \mu\chi + jD^c\chi \\ &= \mu(\chi - j\phi) \\ &= -j\mu(\phi + j\chi). \end{aligned}$$

There is certainly no complex 4×4 matrix j which anticommutes with the complex number i . But let us pass to real 4×4 matrices and identify i_k with its first representation $L(i_k)$ (Chapter 8). We shall write j_k for $R(i_k)$, the second representation of i_k , as in Chapter 8, Exercise 7. Then

$$j_1^2 = j_2^2 = j_3^2 = j_3j_2j_1 = -1, \quad j_k i_l = i_l j_k \quad (k, l = 1, 2, 3).$$

Now replace the complex number i by the real matrix j_1 and write j_2 for j . Then the assumption made above is justified. Putting

$$\psi = \phi + j_2\chi = \phi_0 + j_1\phi_1 + j_2\chi_0 + j_3\chi_1,$$

we may write the above equation as follows:

$$D\psi + j_2\mu\psi = 0.$$

This is essentially Dirac's equation for the electron.

It can be shown that the sixteen matrices $1, i_k, j_l, i_k j_l$ ($k, l = 1, 2, 3$) are linearly independent (e.g., Jacobson [1980] p. 218, Theorem 4.6). Thus ψ is just an arbitrary real 4×4 matrix. However, as we may multiply Dirac's equation by the column vector $(1000)^t$, we may assume, without loss of generality, that ψ itself is a column vector $(\psi_0 \psi_1 \psi_2 \psi_3)^t$ with real components. Nothing prevents us from allowing the ψ_k to be complex numbers, but, as far as the present analysis is concerned, there is no compelling reason for doing so. (However, complex values are forced upon us, as soon as we look at the electron in an electromagnetic field.)

Since a Lorentz transformation sends D onto $p^c D p^t$, we want ψ to be transformed to $p\psi$, hence $D\psi$ to $p^c D\psi$ and $j_2\mu\psi$ to $j_2\mu p\psi = p^c j_2\mu\psi$, thus making the Dirac equation Lorentz invariant.

It is important to note that the biquaternions of norm 1, p and $-p$, while yielding the same Lorentz transformation, both sending x onto pxp^{ct} , induce distinct transformations on ψ , sending it to $p\psi$ and $-p\psi$, respectively. This is the mathematical reason for saying that the electron has *spin* $\frac{1}{2}$.

Exercises

1. If the biquaternion a is viewed as a real 4×4 matrix, show that its transpose is not a^t , but a^{ct} , hence a point in Minkowski space is represented not by a Hermitian, but by a symmetric matrix. A general symmetric matrix has the form

$$x = x_0 + \sum_{k,l=1}^3 x_{kl} i_k j_l,$$

but, for a point in Minkowski space $x_{kl} = 0$ unless $k = l = 1$, since Minkowski space has only four and not ten dimensions. (One version of *string theory* does allow for ten dimensions, presumably for quite different reasons.)

2. Maxwell predicted from his equations that electromagnetic energy is propagated in waves. Einstein used the equation $E = mc^2$ to predict that mass is convertible into energy. Dirac observed that in his equation μ might as well be replaced by $-\mu$ and predicted that the electron must have an anti-particle, now called the positron. Discuss to what extent such predictions from mathematical symbolism to physical reality are justified.

12

Cardinal Numbers

The ideas in this chapter (and Chapter 13) are due to Georg Cantor (1845–1918). Many mathematicians at first rejected Cantor’s work for ideological reasons, claiming that there could be no ‘actual infinity’ in mathematics. Cantor found it impossible to get a decent job and, in 1884, suffered a mental breakdown from which he never fully recovered.

The *cardinal numbers* include the natural numbers $0, 1, 2, \dots$, but go beyond them by including various kinds of infinity. We have some of the same problems in defining a cardinal number as we had with the number 2. Roughly speaking, a cardinal number tells you how many elements there are in a given set. This notion is clear enough for finite sets, but for infinite sets we need some more discussion.

Two sets are said to have the *same cardinality* if and only if there is a one-to-one correspondence between them:

$A \cong B \iff$ there is some $f : A \rightarrow B$ such that f is one-to-one and onto (or *injective* and *surjective*).

For example, $\{3, 5, 7\}$ has the same cardinality as $\{0, 1, 2\}$. Note that \cong is an equivalence relation.

If n is a positive integer, we say that a set A has *cardinality* n if and only if $A \cong \{0, 1, 2, \dots, n-1\}$, and we write $|A| = n$. The cardinality of the empty set is 0: $|\emptyset| = 0$.

A set A is *finite* when $|A|$ is a natural number. Otherwise A is *infinite*. For example, \mathbf{N} is infinite.

We say a set A has cardinality \aleph_0 if and only if $A \cong \mathbf{N}$. Such sets are called ‘countably infinite’ or just ‘countable’. For example, $f(m) = 2m$ is a 1-to-1 function from the set \mathbf{N} of natural numbers onto the set E of even

numbers. Thus the cardinality of E is \aleph_0 and we write $|E| = \aleph_0$. The above example is a special case of what is usually called ‘Galileo’s Theorem’:

If a set S has cardinality \aleph_0 , then any infinite subset of S has cardinality \aleph_0 .

Galileo’s theorem might lead one to think that there is only one infinite cardinal. However, one of Cantor’s achievements was to show that the power set of any set has a higher cardinality than the set itself. By the *power set* of a set S , we mean the set $P(S)$ whose members are the subsets of S .

We write $|A| < |B|$ (A has a lower cardinality than B) if and only if $A \not\cong B$ but B has a proper subset C such that $A \cong C$.

Note that $P(S)$ always has a proper subset C , the set of singletons of S , such that $C \cong S$.

We now have Cantor’s Theorem:

Theorem 12.1. *Any set A has a lower cardinality than $P(A)$.*

Proof: To obtain a contradiction, suppose there is a 1-to-1, onto function $f : A \rightarrow P(A)$. Then every subset of A has the form $f(a)$ for some $a \in A$.

Let $S = \{x \in A \mid x \notin f(x)\}$. This is a subset of A , so there is some $a \in A$ such that $S = f(a)$. If $a \in S$ then, by the defining property of S , $a \notin f(a) = S$, hence $a \notin S$. But then $a \notin f(a)$ and so $a \in S$. Contradiction.

If a finite set has n elements then its power set has 2^n elements. (This is because, in forming a subset, there are two choices for each of the elements in the original set: include it in the subset or leave it out.) In general, if a set has cardinality k , we use the symbol 2^k to denote the cardinality of its power set. For example, $|P(\mathbf{N})| = 2^{\aleph_0}$. By Cantor’s Theorem,

$$\aleph_0 < 2^{\aleph_0} < 2^{2^{\aleph_0}} < \dots$$

In other words, there are an infinite number of infinities. (It was statements like these that got Cantor into trouble with Kronecker.)

Is there some subset of $P(N)$ whose cardinality is greater than \aleph_0 but not as great as 2^{\aleph_0} ? The hypothesis that the answer to this question is NO is called the *continuum hypothesis*. In 1940 Kurt Gödel showed that the continuum hypothesis is consistent with the usual axioms of set theory. In 1963 Paul Cohen showed that the negation of the continuum hypothesis is also consistent with the usual axioms of set theory. In other words, our basic notions about sets neither imply nor preclude the continuum hypothesis. We say that the continuum hypothesis is *independent* of the axioms of set theory.

Exercises

1. Suppose that $f : A \rightarrow B$, $g : B \rightarrow A$, $fg = 1_B$ and $gf = 1_A$ (where 1_S is the identity function on S). Prove that f is one-one and onto. Conversely, if f is one-one and onto, show that it has an inverse g .
2. Let S be the set of natural numbers of the form $2^a 3^b$ where a and b are positive integers whose gcd is 1. Use Galileo's theorem to show that $|S| = \aleph_0$.
3. Using Exercise 2, show that the set of positive rationals has cardinality \aleph_0 .
4. Let T be the set of all sequences formed from 0's and 1's. Show that $|T| = 2^{\aleph_0}$.
5. Let T be the set of reals from 0 to 1. Show that $|T| = 2^{\aleph_0}$.
6. Prove Galileo's Theorem.
7. Prove that the set of all functions from \mathbf{N} to \mathbf{N} has higher cardinality than \mathbf{N} .

13

Cardinal Arithmetic

We begin by defining three binary operations for sets in general:

$$A \times B = \{(a, b) | a \in A \text{ and } b \in B\},$$

$$A^B = \{f | f : B \rightarrow A\},$$

$$A + B = (A \times 0) \cup (B \times 1).$$

Note that $A \times \emptyset = \emptyset$, and, if $B \neq \emptyset$, $\emptyset^B = \emptyset$. These definitions are motivated by the fact that, for finite sets A and B ,

$$|A \times B| = |A| \times |B|,$$

$$|A^B| = |A|^{|B|},$$

$$|A + B| = |A| + |B|.$$

We also have the following theorem, which generalizes the results of Chapter 2. We have written 0 for \emptyset and 1 for $\{\emptyset\}$.

Theorem 13.1.

1. $A + B \cong B + A$,
2. $(A + B) + C \cong A + (B + C)$,
3. $A \times B \cong B \times A$,
4. $(A \times B) \times C \cong A \times (B \times C)$,

5. $A \times (B + C) \cong (A \times B) + (A \times C)$,
6. $(A^B)^C \cong A^{C \times B}$,
7. $A^B \times A^C \cong A^{B+C}$,
8. $(A \times B)^C \cong A^C \times B^C$,
9. $A + 0 \cong A$,
10. $A \times 1 \cong A$,
11. $A \times 0 \cong 0$,
12. $A^0 \cong 1$, $A^1 \cong A$ and, if $B \neq 0$, then $0^B \cong 0$.

We shall not give a complete proof of this theorem, but, as an example, we shall prove that $(A^B)^C \cong A^{C \times B}$.

Let f be any member of $(A^B)^C$, that is, let f be any function such that $f : C \rightarrow A^B$. Then, if c is a typical element of C , $f(c) \in A^B$. Hence, if b is a typical element of B , $(f(c))(b) \in A$. Let g be any member of $A^{C \times B}$, that is, let g be any function such that $g : C \times B \rightarrow A$. Then, if (c, b) is a typical element of $C \times B$ (with $c \in C$, $b \in B$), $g((c, b)) \in A$.

We define $F : A^{C \times B} \rightarrow (A^B)^C$ such that $((F(g))(c))(b) = g((c, b))$. We define $G : (A^B)^C \rightarrow A^{C \times B}$ such that $(G(f))((c, b)) = (f(c))(b)$.

Then FG is just the identity function on $(A^B)^C$. For $(FG)(f) = F(G(f))$ and this equals f just in case, for any $c \in C$, $(F(G(f)))(c) = f(c)$. Moreover, the last equation is true just in case, for any $b \in B$, $((F(G(f)))(c))(b) = (f(c))(b)$. Now $((F(G(f)))(c))(b) = (G(f))((c, b))$, by the definition of F . Furthermore, by the definition of G , $(G(f))((c, b)) = (f(c))(b)$ as required.

Again, GF is just the identity function on $A^{C \times B}$. For $(GF)(g) = G(F(g))$ and this equals g just in case, for any $(c, b) \in C \times B$, $(G(F(g)))(c, b) = g((c, b))$. By the definition of G , $(G(F(g)))(c, b) = ((F(g))(c))(b)$, and, by the definition of F this equals $g((c, b))$ as required.

Hence F is one-one and onto (Chapter 12, Exercise 1). Applying the function F is called *currying* in Computer Science, after the logician Haskell B. Curry.

Here are some hints for proving the remaining 11 parts of the above theorem, skipping a few unnecessary parentheses.

(1) An element of the left-hand side must be of the form $(a, 0)$ or $(b, 1)$, where $a \in A$ and $b \in B$. Let $F(a, 0) = (a, 1)$ and $F(b, 1) = (b, 0)$. Find the inverse G of F .

(3) An element of the left-hand side has the form (a, b) , where $a \in A$ and $b \in B$. Let $F(a, b) = (b, a)$ and find the inverse G of F .

(5) An element of the left-hand side has the form $(a, (b, 0))$ or $(a, (c, 1))$, where $a \in A$, $b \in B$ and $c \in C$. Let $F(a, (b, 0)) = ((a, b), 0)$, $F(a, (c, 1)) = ((a, c), 1)$ and find the inverse G of F .

(7) An element of the left-hand side has the form (f, g) , where $f : B \rightarrow A$ and $g : C \rightarrow A$. Define $F(f, g)$ by stipulating $F(f, g)(b, 0) = f(b)$ and $F(f, g)(c, 1) = g(c)$. An element of the right-hand side has the form $h : B + C \rightarrow A$. Define $G(h)$ as the pair $(G(h)_0, G(h)_1)$, where $G(h)_0(b) = h(b, 0)$ and $G(h)_1(c) = h(c, 1)$. Show that G is the inverse of F .

(8) An element of the left-hand side has the form $h : C \rightarrow A \times B$. From this we obtain two functions $h_0 : C \rightarrow A$ and $h_1 : C \rightarrow B$ such that $h(c) = (h_0(c), h_1(c))$. Let $F(h) = (h_0, h_1)$. An element of the right-hand side is a pair (f, g) , where $f : C \rightarrow A$ and $g : C \rightarrow B$. Let $G(f, g)(c) = (f(c), g(c))$ and show that G is the inverse of F .

(12) An element of A^B is a function $f : B \rightarrow A$, that is, a subset f of $A \times B$ such that, for every $b \in B$, there exists a unique $a \in A$ such that $(a, b) \in f$. Show that there is exactly one function $\emptyset \rightarrow A$, that the functions $\{\emptyset\} \rightarrow A$ are in one-to-one correspondence with the elements of A , and, finally that, when $B \neq \emptyset$, there is no function $B \rightarrow \emptyset$.

We can now define equivalent binary operations for cardinals, so that

$$\begin{aligned} |A| + |B| &= |A + B|, \\ |A| \times |B| &= |A \times B|, \\ |A|^{|B|} &= |A^B|. \end{aligned}$$

For example, $\aleph_0 + 2 = \aleph_0$ since $\mathbf{N} + \{0, 1\} \cong \mathbf{N}$.

Exercises

1. Prove the remaining eleven parts of the above theorem.
2. Prove that $|A| \times (|B| + |C|) = (|A| \times |B|) + (|A| \times |C|)$.
3. Prove that $\aleph_0 + \aleph_0 = \aleph_0$.
4. Prove that $\aleph_0 \times 2 = \aleph_0$.
5. Prove that $\mathbf{N} + \mathbf{N}$ and $\mathbf{N} \times \mathbf{N}$ have the same cardinality as \mathbf{N} .
6. Prove that $\aleph_0^2 = \aleph_0$.
7. Prove that $|A| \times |A| \times |A| = |A|^3$.
8. Simplify $2^{\aleph_0} \times 2^{\aleph_0}$.
9. Simplify 8^{\aleph_0} .

10. Prove that every infinite set has a subset which can be placed in one-one correspondence with \mathbf{N} .
11. Project: Look up the Schroeder-Bernstein Theorem and prove:

$$\aleph_0^{\aleph_0} = 2^{2^{\aleph_0}}.$$

Continued Fractions

As Fowler shows in *The Mathematics of Plato's Academy*, continued fractions are implicit in ancient Greek mathematics. As far as we know, they were not explicitly defined before 1618, when Daniel Schwenter rediscovered them.

Continued fractions are implicit in Euclid's Algorithm (Book VII, Proposition 2) dating from 300 BC. This is a procedure for finding the greatest common divisor (gcd) of two positive integers. The *greatest common divisor* of two positive integers a and b is the positive integer d whose divisors are precisely the common divisors of a and b . In other words, d divides a and b and any common divisor of a and b divides d . This notion can easily be extended to natural numbers or even integers, but one has to be careful. For example, $\gcd(0, 17) = 17$, but $\gcd(0, 0) = 0$, even though 17 is also a common divisor of 0 and 0, yet $17 > 0$. For the gcd of two nonzero integers, say 12 and -15 , one has two candidates that meet the above definition, 3 and -3 in this case, and one usually chooses the positive one. The algorithm is best described with the help of an example. Suppose we want to find the gcd of 502 and 1604. We perform four divisions as follows.

	3		5		8		6
502	1604	98	502	12	98	2	12
	1506		490		96		12
	98		12		2		0

Since $\text{dividend} = (\text{divisor} \times \text{quotient}) + \text{remainder}$, any number which divides dividend and divisor will also divide the remainder. Thus any factors common to the two original integers are carried through the calculations

until they surface in the last nonzero remainder. The last nonzero remainder is thus the gcd of the original integers. The above calculations can be rewritten as follows:

$$\frac{1604}{502} = 3 + \frac{1}{\frac{502}{98}} = 3 + \frac{1}{5 + \frac{12}{98}} = 3 + \frac{1}{5 + \frac{1}{\frac{98}{12}}} = 3 + \frac{1}{5 + \frac{1}{8 + \frac{2}{12}}} = 3 + \frac{1}{5 + \frac{1}{8 + \frac{1}{6}}}$$

The final expression is a (*simple*) *continued fraction*.

Standard abbreviations for continued fractions allow us to write the above fraction as $3 + \frac{1}{5 + \frac{1}{8 + \frac{1}{6}}}$ or as $(3, 5, 8, 6)$. Thus (a_0, a_1, a_2) is the fraction

$$a_0 + \frac{1}{a_1 + \frac{1}{a_2}}.$$

When one divides by a positive integer, one always obtains a remainder which is less than the divisor. Thus the sequence of divisors in Euclid's Algorithm steadily decreases. Using Euclid's Algorithm, we can write any positive rational number as a finite continued fraction $(a_0, a_1, a_2, \dots, a_n)$, where the a_i are natural numbers and, for $i > 0$, $a_i > 0$.

However, even irrational real numbers can be written as continued fractions, and this too was known to the ancient Greeks, according to Fowler. We then obtain infinite continued fractions (a_0, a_1, a_2, \dots) ; but these may be approximated by finite continued fractions $c_0 = a_0$, $c_1 = (a_0, a_1)$, $c_2 = (a_0, a_1, a_2)$, etc., called *convergents*. Ultimately we are interested in the case when the a_n are all positive integers, with the possible exception of a_0 . For the following argument only, we shall allow the a_n to be positive rational numbers for $n > 0$.

It is easily seen that $c_0 = \frac{a_0}{1}$ and $c_1 = \frac{a_0 a_1 + 1}{a_1}$. To calculate c_n when $n > 1$ we define two sequences of rationals (ultimately integers):

$$p_0 = a_0, p_1 = a_0 a_1 + 1, p_n = a_n p_{n-1} + p_{n-2} \text{ if } n > 1,$$

$$q_0 = 1, q_1 = a_1, q_n = a_n q_{n-1} + q_{n-2} \text{ if } n > 1.$$

We claim that $c_n = \frac{p_n}{q_n}$ for all n . This is evidently so when $n = 0$ or 1 , so assume the result for n for any admissible (a_0, a_1, \dots, a_n) ; we shall show that it also holds for $n + 1$. Now c_{n+1} can be obtained from

$$c_n = \frac{a_n p_{n-1} + p_{n-2}}{a_n q_{n-1} + q_{n-2}}$$

by replacing a_n by $a_n + \frac{1}{a_{n+1}}$, with the help of the induction assumption applied to $(a_0, \dots, a_n + \frac{1}{a_{n+1}})$. Multiplying top and bottom of the resulting

ratio by a_{n+1} , the top becomes

$$\begin{aligned} a_{n+1} \left(\left(a_n + \frac{1}{a_{n+1}} \right) p_{n-1} + p_{n-2} \right) &= a_{n+1}(a_n p_{n-1} + p_{n-2}) + p_{n-1} \\ &= a_{n+1} p_n + p_{n-1} = p_{n+1}. \end{aligned}$$

Similarly, the bottom becomes q_{n+1} , so $c_{n+1} = p_{n+1}/q_{n+1}$, and therefore the result holds by mathematical induction.

The proof involved the rational number $a_n + \frac{1}{a_{n+1}}$ in the induction assumption, but from now on we shall stick to integer a_n . We have proved:

Theorem 14.1. *If $a_0 \in \mathbf{N}$ and $0 < a_n \in \mathbf{N}$ for $n > 0$, the n th convergent of (a_0, a_1, a_2, \dots) is p_n/q_n , where the p_n and q_n are defined inductively as above.*

Note that, if the a_n are all natural numbers, as we are now supposing, the inductively defined p_n and q_n will be positive integers, in fact, strictly increasing sequences of positive integers.

Theorem 14.2. *Let p_n and q_n be defined inductively as above. For all $n > 0$, $p_n q_{n-1} - p_{n-1} q_n = (-1)^{n-1}$.*

Proof: This is evidently so when $n = 1$ or 2 . Assume the result for n . Then

$$\begin{aligned} p_{n+1} q_n - p_n q_{n+1} &= (a_{n+1} p_n + p_{n-1}) q_n - p_n (a_{n+1} q_n + q_{n-1}) \\ &= p_{n-1} q_n - p_n q_{n-1} \\ &= -(-1)^{n-1} \\ &= (-1)^n. \end{aligned}$$

An immediate consequence of this theorem is

Corollary 14.3. *For all $n \in \mathbf{N}$, $\gcd(p_n, q_n) = 1$; in other words, p_n/q_n is in lowest terms.*

Another immediate consequence is

Corollary 14.4. *For all $n > 0$, $c_n - c_{n-1} = \frac{p_n}{q_n} - \frac{p_{n-1}}{q_{n-1}} = \frac{(-1)^{n-1}}{q_n q_{n-1}}$.*

It follows that, for odd n ,

$$c_n - c_{n-2} = \frac{1}{q_n q_{n-1}} - \frac{1}{q_{n-2} q_{n-1}} < 0,$$

since the q_n are strictly increasing, and similarly, for even n , $c_n - c_{n-2} > 0$. Thus we have a strictly increasing sequence

$$c_1 > c_3 > c_5 > \dots$$

and a strictly decreasing sequence

$$c_0 < c_2 < c_4 < \cdots.$$

Moreover, by Corollary 14.4, the difference of the two sequences tends to 0, hence they must have a common limit a and we write $a = (a_0, a_1, a_2, \dots)$.

For example, the infinite continued fraction $(1, 1, 1, \dots)$ has a limit x , hence $x = 1 + \frac{1}{x}$, so $x^2 - x - 1 = 0$, hence $x = \frac{1}{2}(1 + \sqrt{5})$.

Exercises

1. Find the gcd of 10403 and 2987.
2. Show that $\gcd(a, b) = \gcd(a, a - b)$.
3. Calculate $(1, 2, 1, 2, 1, 2, \dots)$.
4. Show that every positive rational number b/c can be written as a simple continued fraction with an even number of a_i .
5. Let b/c be a positive proper fraction in lowest terms. Using Exercise 4, show that b/c can be written in the form $1/d + e/f$ where $b > e$ and $d > c > f$.
6. Using Exercise 5, show that every proper reduced fraction can be expressed as a sum of distinct *unit fractions* (that is, fractions with numerator 1 and a positive integer denominator).
7. Use the method of Exercises 5 and 6 to express $67/120$ as a sum of distinct unit fractions.
8. Associate with each letter the number of its place in the (English) alphabet. Then each word is associated with a sequence a_0, a_1, \dots, a_n . Encode the word into (a_0, a_1, \dots, a_n) , simplified into an ordinary fraction. How would you decipher such a fraction?
9. Show that $(2a, a, 2a, a, 2a, a, \dots) = a + \sqrt{a^2 + 2}$.

The Fundamental Theorem of Arithmetic

One can use Euclid's algorithm to find the gcd of two positive integers a and b . One can also exploit the algorithm to express the gcd d in the form $d = ax + by$ where x and y are integers. Actually, d is the smallest positive integer with this property, and this fact can also be used to describe the gcd of a and b .

Theorem 15.1. *Given positive integers a and b , their gcd is the smallest positive integer d such that $d = ax + by$ with $x, y \in \mathbf{Z}$.*

Proof: Note that the set $\{ax + by | x, y \in \mathbf{Z}\}$ does contain positive integers, e.g., $2ab$. Let d be the smallest positive integer in the set. Clearly any common divisor of a and b divides d . So, to prove that $d = \gcd(a, b)$, we only have to show that d divides a and b . To prove that d divides a , divide a by d to get quotient q and remainder r , so that

$$a = qd + r, \quad 0 \leq r < d.$$

Then $r = a - qd = a - q(ax + by) = a(1 - qx) + b(-qy)$ is also of the form $ax' + by'$ with $x', y' \in \mathbf{Z}$. Since d was the smallest positive integer of this form and since $0 \leq r < d$, it follows that $r = 0$, hence d divides a . Similarly, d divides b , so d is a common divisor of a and b .

The following consequence of the theorem is known as the *Fundamental Lemma of Arithmetic*.

Lemma 15.2. *Given positive integers a , b and c , if a divides bc and*

$\gcd(a, b) = 1$, then a divides c . In particular, if a prime number p divides bc , then p divides b or p divides c .

Proof: If $\gcd(a, b) = 1$, it follows from the theorem that there are integers x and y such that $ax + by = 1$. Moreover, if a divides bc , there is an integer z such that $bc = az$. Therefore,

$$c = axc + byc = a(xc + yz),$$

and so a divides c .

The Fundamental Theorem of Arithmetic asserts that every positive integer has a factorization into primes; moreover, if we disregard the order of the primes, this factorization is unique. This theorem follows quite easily from the above lemma and we leave it as an exercise. For another proof and a discussion of its history, see Part I, Chapter 3.

Given positive integers a and b , say with $a > b$, we obtain a continued fraction expansion $a/b = (a_0, a_1, \dots, a_n)$. Even though this continued fraction is *finite*, we can still use the analysis of the previous section to calculate its convergents p_0/q_0 up to p_n/q_n . In particular, $a/b = p_n/q_n$, hence $aq_n = bp_n$. Since $\gcd(p_n, q_n) = 1$, it follows from the Fundamental Lemma of Arithmetic that p_n divides a , say $a = p_nd$, whence also $b = q_nd$. Evidently, d is the gcd of a and b . But it also follows from Theorem 14.2, upon multiplying by d , that $aq_{n-1} - bp_{n-1} = d(-1)^{n-1}$. This allows us to find particular integers x and y such that $d = ax + by$, namely,

$$x = (-1)^{n-1}q_{n-1}, \quad y = (-1)^np_{n-1}.$$

Exercises

1. Prove that the smallest divisor > 1 of a positive integer > 1 is a prime number.
2. Deduce that every positive integer > 1 is a product of prime numbers and use the Fundamental Lemma of Arithmetic to show that this factorization into primes is unique.
3. If a, b and c are positive integers and $c > 1$, show that

$$\gcd(c^a - 1, c^b - 1) = c^d - 1,$$

where $d = \gcd(a, b)$. (Hint: use the above theorem.)

Linear Diophantine Equations

A *linear Diophantine equation* in two variables is an equation of the form $ax + by = c$ where a, b and c are given integers, and x and y are unknown integers. Sometimes x and y are restricted to the set of *positive* integers.

For example, $4x + 6y = 8$ is a linear Diophantine equation. It has solution $x = 2$ and $y = 0$ (among others). Here we are not interested in noninteger solutions such as $x = \frac{1}{2}$, $y = 1$.

In order to simplify the presentation, we shall allow negative numbers as solutions, but we shall take a, b and c to be positive. By dividing the material into different cases, we could elaborate the whole theory in terms of positive integers. Thus, there is nothing in this chapter which would be inaccessible to the ancient Greeks. On the contrary, it is probable that they used essentially the following method to attack these equations.

The adjective ‘Diophantine’ comes from the name ‘Diophantus’. Diophantus of Alexandria lived about 250 AD. However, in his equations he did not restrict x and y to be integers, but allowed them to be rationals. It was Brahmagupta of India (628 AD) who gave the first complete solution to the linear Diophantine equation (Boyer [1989], pp. 244-47).

Let $d = \gcd(a, b)$. Then there are integers a' and b' such that $a = da'$ and $b = db'$. If $ax + by = c$ then $d(a'x + b'y) = c$ and hence d is a factor of c . Thus, if d is not a factor of c , the Diophantine equation has no solutions. On the other hand, if d is a factor of c , then $c = dc'$ for some integer c' , and we can cancel d to get $a'x + b'y = c'$, with $\gcd(a', b') = 1$. In giving a solution of $ax + by = c$, we can thus, without loss of generality, begin by assuming that $\gcd(a, b) = 1$.

To solve the Diophantine equation $ax + by = c$ when $\gcd(a, b) = 1$:

1. Use Euclid's algorithm to write $a/b = (a_0, a_1, \dots, a_n)$.
2. Since a/b is in lowest terms, $a = \pm p_n$ and $b = \pm q_n$. Since $p_n q_{n-1} - p_{n-1} q_n = \pm 1$, it follows that, with some selection of signs, $x_0 = \pm c q_{n-1}$, $y_0 = \pm c p_{n-1}$ is a solution of the Diophantine equation. Thus, using the recursive definitions given in Chapter 14, calculate p_{n-1} and q_{n-1} and pick the signs so that x_0, y_0 is a solution.
3. If t is any integer, $x = x_0 + bt$, $y = y_0 - at$ is another solution of the Diophantine equation.
4. There are no integer solutions other than those mentioned in (3) above. For if $ax + by = c = ax_0 + by_0$, then $a(x - x_0) = b(y_0 - y)$ and b is a factor of $x - x_0$, since $\gcd(a, b) = 1$. Hence, for some integer t , $x = x_0 + bt$.

Example: Solve $25x + 55y = 50$.

First we get the gcd of 25 and 55 by using Euclid's algorithm:

$$25/55 = 0 + \frac{1}{55/25} = 0 + \frac{1}{2 + 5/25} = 0 + \frac{1}{2 + 1/5} = (0, 2, 5).$$

The gcd is the last divisor, namely, 5. (It is merely a coincidence that the last partial quotient is also 5.) We note that 5 divides 50 and so the equation does have integer solutions. We factor out the 5, obtaining $5x + 11y = 10$. We calculate the penultimate convergent $p_1/q_1 = (0, 2) = 1/2$. With the right signs, one solution is $x_0 = \pm 10 \times 2$ and $y_0 = \pm 10 \times 1$. Indeed, we can take $x_0 = -20$ and $y_0 = 10$. The general solution is $x = -20 + 11t$, $y = 10 - 5t$, where t is any integer.

Sometimes we want only positive solutions. In that case we must have $x_0 + bt > 0$ and $y_0 - at > 0$. In the above example this would require an integer t such that $1 < 20/11 < t < 10/5 = 2$, which is impossible.

Exercises

1. Solve the Diophantine equation $101x + 753y = 100,000$. (There are two positive integer solutions.)
2. Solve the Diophantine equation $158x + 57y = 20,000$. (There are two positive integer solutions.)
3. Solve the Diophantine equation $91x + 221y = 1053$. (There are no positive integer solutions.)

4. Show that the following Diophantine equation has a unique solution in positive integers: $17x + 19y = 320$.
5. The Sultana used to divide her maids into two companies, one which would follow her five abreast and the other which would follow her seven abreast – both in rectangular formation. These companies would consist of different numbers of maids on each of nine different days. What is the smallest number of maids the Sultana could have had?
6. Show that, if a and b are positive integers, then $ax + by = c$ has no solutions in positive integers when $[-x_0/b] \geq [y_0/a]$. (Here x_0 and y_0 are the solutions described above, and $[z]$ is the greatest integer not exceeding z .)

Quadratic Surds

Let d be a positive nonsquare integer. Using the Fundamental Theorem of Arithmetic, it is not hard to show that \sqrt{d} is irrational. This was first proved by Theaetetus in about 400 BC.

If a and $b \neq 0$ are integers, the expression $(a + \sqrt{d})/b$ is called a *quadratic surd*. Note that if a' and b' are integers, and $(a + \sqrt{d})/b = (a' + \sqrt{d})/b'$ then $a = a'$ and $b = b'$. The proof is left as an exercise.

In this chapter we study continued fraction expansions of quadratic surds. These expansions can be used to solve quadratic Diophantine equations, such as the ‘Pell equation’ $x^2 - dy^2 = 1$.

To begin with an example, suppose that $x = (1, 1, 1, \dots)$. Then $x = 1 + 1/x$ and hence $x^2 - x - 1 = 0$, with the result that $x = \frac{1}{2}(1 + \sqrt{5})$. (We cannot take the other root of the equation since x is positive.) This is a quadratic surd. Indeed, it is a very famous one known as ‘the golden ratio’. It is the ratio of the side to the base in the triangles obtained by connecting the five points of the Pythagorean star.

As another example, suppose that $y = (1, 1, 2, 1, 2, 1, \dots)$. To evaluate this continued fraction, let $x = (1, 2, 1, 2, 1, \dots)$. Then $y = 1 + 1/x$. Furthermore, $x = 1 + \frac{1}{2+1/x}$, so that $2x^2 - 2x - 1 = 0$ and hence $x = \frac{1}{2}(1 + \sqrt{3})$. Thus $y = \sqrt{3}$.

Generalizing from these two examples, it is easily shown that any continued fraction which is ultimately periodic represents a quadratic surd. Less obvious is the converse of this statement, namely, that every quadratic surd has a continued fraction expansion which is ultimately periodic. This was first proved by Lagrange in about 1770. We shall prove the special case:

Theorem 17.1. *If d is a positive integer which is not a perfect square, the continued fraction expansion of \sqrt{d} is ultimately periodic.*

Proof: We define sequences of integers a_n and b_n and a sequence of rationals r_n as follows, where $[\rho]$ denotes the greatest integer in ρ :

$$a_0 = 0, r_0 = 1, b_n = \left[\frac{\sqrt{d} + a_n}{r_n} \right], a_{n+1} = b_n r_n - a_n, r_{n+1} = (d_n - a_{n+1}^2)/r_n.$$

Then it is easily verified that

$$\frac{\sqrt{d} + a_n}{r_n} = b_n + \frac{1}{b_{n+1} + \frac{1}{b_{n+2} + \dots}},$$

in particular,

$$\sqrt{d} = b_0 + \frac{1}{b_1 + \frac{1}{b_2 + \dots}},$$

and hence $b_{n+1} > 0$. It also follows by mathematical induction that the r_n are integers, once it is realized that

$$r_{n+2} = r_n - b_{n+1}^2 r_{n+1} + 2b_{n+1} a_{n+1}.$$

Let $x_n = (\sqrt{d} + a_n)/r_n$ and let t_n be its conjugate $(-\sqrt{d} + a_n)/r_n$. From Chapter 14 we know that

$$\sqrt{d} = \frac{x_n p_{n-1} + p_{n-2}}{x_n q_{n-1} + q_{n-2}}.$$

Taking conjugates, we obtain

$$-\sqrt{d} = \frac{t_n p_{n-1} + p_{n-2}}{t_n q_{n-1} + q_{n-2}}.$$

Solving for t_n , we obtain

$$t_n = \left(-\frac{q_{n-2}}{q_{n-1}} \right) \left(\frac{\sqrt{d} + p_{n-2}/q_{n-2}}{\sqrt{d} + p_{n-1}/q_{n-1}} \right).$$

As n increases, the second factor tends to 1. The q 's are positive, so, for sufficiently large n , $t_n < 0$. Now $x_n > 1$ so, for sufficiently large n ,

$$\frac{2\sqrt{d}}{r_n} = x_n - t_n > 1$$

and hence $r_n > 0$ and $2\sqrt{d} > r_n$. Since $t_n < 0$, it then follows that $a_n < \sqrt{d}$ while $x_n > 1$ implies $a_n > -\sqrt{d}$. Since

$$2\sqrt{d} > r_n > 0,$$

$$\sqrt{d} > a_n > -\sqrt{d},$$

there are only a finite number of possibilities for $x_n = (\sqrt{d} + a_n)/r_n$ and so the simple continued fraction must repeat.

Exercises

1. Find the continued fraction expansions of $\sqrt{2}$ and $\sqrt{10}$.
2. Find the quadratic surd which is represented by the periodic continued fraction $(a, b, c, a, b, c, \dots)$.
3. Fill in the details in the proof of the above theorem.
4. Show that for nonsquare d , the continued fraction expansion of \sqrt{d} is of the form $\sqrt{d} = (a_0, \overline{a_1, \dots, a_{n-1}, 2a_0})$, with the part under the bar periodic.
5. For d as in Exercise 4, show that $p_{n-1}^2 - dq_{n-1}^2 = \pm 1$, where p_{n-1}/q_{n-1} is the $(n-1)$ th convergent of \sqrt{d} . (Hint: show that

$$\sqrt{d} = \frac{\alpha_n p_{n-1} - p_{n-2}}{\alpha_n q_{n-1} - q_{n-2}}$$

where $\alpha_n = \sqrt{d} + a_0$.)

6. Find a positive integer solution to the equation $x^2 - 61y^2 = \pm 1$.
7. Prove that in $\frac{a+\sqrt{d}}{b}$, the positive integers a and b are uniquely determined.

Note that Exercises 4 and 5 are more difficult.

Pythagorean Triangles and Fermat's Last Theorem

A *Pythagorean triangle* is a right triangle all of whose sides have integer lengths. For example, let ABC be a triangle with a right angle at vertex C . Suppose that the sides AC and BC have lengths 5 and 12, respectively. Then the hypotenuse has length $\sqrt{5^2 + 12^2} = 13$. Since all three sides have integer lengths, this is a Pythagorean triangle.

Note that if k is any positive integer, the triangle with sides of lengths $5k$, $12k$ and $13k$ is a Pythagorean triangle too. This triangle is just a magnification of the previous one and so it is not very interesting. However, it does suggest that it would be worthwhile to find all the Pythagorean triangles whose gcd is 1. These are called *primitive* Pythagorean triangles. It is easy to show that if x , y and z are the sides of a right triangle, then $\gcd(x, y, z) = \gcd(x, y) = \gcd(y, z)$. In what follows we shall assume that this gcd is equal to 1.

We shall use the following result:

Lemma 18.1. *If m and n are nonnegative integers such that $\gcd(m, n) = 1$ and mn is a square, then both m and n are squares.*

Professor Tournesol discovered an amusing proof of this theorem while he was in jail for having failed the president's son. The prisoners were put in a long row of cells. At first all the doors were unlocked, but then the jailor walked by and locked every second door. He walked by again and stopped at every third door, locking it if it was unlocked, but unlocking it if it was locked. On his next round he stopped at every fourth door, locking it if it was unlocked, unlocking it if it was locked, and so on. Professor Tournesol soon realised that the m th cell would be unlocked in the end just in case

m had an odd number of divisors. Now, if d divides m then so does m/d and it would seem that the divisors of m come in pairs. Unless... 'what if $d = m/d$?' thought the professor, 'then the divisor d does not pair off with another, and $d = m/d$ just in case m is a square.'

Let $\tau(x)$ be the number of divisors of a positive integer x . Then $\tau(x)$ is odd just in case x is a square. Now, if $\gcd(m, n) = 1$ then a typical divisor of mn is of the form dg where d divides m and g divides n . Thus $\tau(mn) = \tau(m)\tau(n)$. If mn is a square then $\tau(mn)$ is odd, and hence both $\tau(m)$ and $\tau(n)$ are odd. Thus m and n are both squares. QED.

In order to solve the Diophantine equation $x^2 + y^2 = z^2$ with $\gcd(x, y, z) = 1$, first note that x and y are not both odd. For if $x = 2a + 1$ and $y = 2b + 1$ then $x^2 + y^2 = 4(a^2 + b^2 + a + b) + 2$ which is not a square, since squares of odd numbers have the form $4(c^2 + c) + 1$ and squares of even numbers have the form $4c^2$. Since $\gcd(x, y) = 1$, it follows that x and y are not both even either. Hence exactly one of x and y is even. Without loss of generality, let us say that $y = 2y'$ and that x is odd. Then x^2 is also odd and y^2 is even. It follows that $z^2 = x^2 + y^2$ is odd and thus z is odd. Since x and y are both odd, $\frac{1}{2}(z + x)$ and $\frac{1}{2}(z - x)$ are both integers. Moreover, their gcd is 1, since any factor which divides them both also divides their sum z and their difference x . But $\gcd(x, z) = 1$.

Since $x^2 + y^2 = z^2$, we have $\frac{1}{2}(z + x)\frac{1}{2}(z - x) = y'^2$. From Professor Tournesol's discovery it follows that there are positive integers u and v with $\gcd(u, v) = 1$ such that $\frac{1}{2}(z + x) = u^2$ and $\frac{1}{2}(z - x) = v^2$. This gives $z = u^2 + v^2$, $x = u^2 - v^2$ and $y = 2y' = 2uv$. Note that u and v are not both odd, since z is not even. The above may be summarized as follows:

Theorem 18.2. *Let x , y and z be positive integers. Then $x^2 + y^2 = z^2$ with y even and $\gcd(x, y, z) = 1$ if and only if, for some positive integers u and v , not both of which are odd, with $u > v$ and $\gcd(u, v) = 1$,*

$$x = u^2 - v^2, \quad y = 2uv, \quad \text{and} \quad z = u^2 + v^2.$$

The sufficiency of the above condition (that is, the 'if' part of the theorem) was known to the Mesopotamians about 4000 years ago (Neugebauer).

The eighth problem in the second book of the *Arithmetica* of Diophantus is to express 16 as a sum of two rational squares. Fermat (1601–1665) had a copy of this book and he enjoyed writing notes in its margins. Unlike Diophantus, Fermat was only interested in (positive) integer solutions to the equations in the *Arithmetica*. Fermat knew that the square of an integer can often be expressed as a sum of two positive integer squares. In the margin beside the problem about expressing 16 as a sum of two squares, Fermat wrote:

On the other hand it is impossible to separate a cube into two cubes, or a biquadratic into two biquadratics, or generally *any power except a square into two powers with the same exponent*. I have discovered a truly marvellous proof of this, which, however, the margin is not large enough to contain' (p. 145, Heath's translation of the *Arithmetica*).

Fermat's assertion is called his 'Last Theorem', although, until quite recently, it would have been safer to call it a 'conjecture'. It is now believed that Fermat only proved the special case when the power $n = 4$, and the 'theorem' remained an open problem for 350 years, though many special cases were proved in that period. The complete theorem was finally proved by Andrew Wiles of Princeton University in 1994. His proof depended on the work of many other mathematicians, notably on a crucial result by K. A. Ribet, as well as ideas from G. Y. Taniyama, G. Shimura, B. Mazur and G. Frey, among others, and the last minute collaboration of Richard Taylor.

Fermat himself showed that $z^4 - x^4 = w^2$ has no solution in positive integers, and from this it follows at once that $x^4 + y^4 = z^4$ has no solution in positive integers (*Oeuvres* I, p. 340 and *Arithmetica*, 2nd edn., p. 293). We shall give a proof of this, due to Euler, which uses Fermat's *method of descent* (from a larger to a smaller solution):

Theorem 18.3. *There are no positive integers x, y and z such that*

$$x^4 + y^4 = z^2.$$

Proof: To obtain a contradiction, suppose there are such integers. Let us take such a triple with the product xy minimized. Then $\gcd(x, y) = 1$. Since x^2 , y^2 and z are the sides of a primitive Pythagorean triangle, exactly one of x and y is even. Without loss of generality, let us say that x is even. By our previous theorem, there are positive integers u and v , not both odd, with $\gcd(u, v) = 1$, such that $x^2 = 2uv$ and $y^2 = u^2 - v^2$. Since $v^2 + y^2 = u^2$ and y is odd, v must be even. Since $\gcd(2v, u) = 1$ and $2vu = x^2$, it follows that $2v$ and u are squares (recall Professor Tournesol). Thus $u = c^2$ for some positive integer c .

Again by our above theorem, there are positive integers s and t , not both odd, with $\gcd(s, t) = 1$, such that $v = 2st$ and $u = s^2 + t^2$. Since $2v$ is a square, so is $2v/4 = v/2 = st$. Thus there are positive integers a and b such that $s = a^2$ and $t = b^2$.

The fact that $u = s^2 + t^2$ implies that $a^4 + b^4 = c^2$. Moreover, $(ab)^2 = st = v/2 < 2uv = x^2 \leq (xy)^2$ so that $ab < xy$. But this contradicts the minimality of xy .

Exercises

1. Prove that if a triangle has sides of lengths x , y and z , and $x^2 + y^2 = z^2$ then the triangle has a right angle.
2. Show that if x , y and z are integers such that $x^2 + y^2 = z^2$ then $\gcd(x, y, z) = \gcd(x, z)$.
3. How many primitive Pythagorean triangles are there with hypotenuse < 50 ?
4. Show that if m and n are positive integers with $\gcd(m, n) = 1$ and mn is a cube then m is a cube.
5. Show that if u , v , u' and v' are positive integers such that $u^2 - v^2 = u'^2 - v'^2$ and $2uv = 2u'v'$ then $u = u'$ and $v = v'$.
6. Solve the Diophantine equation $x^{28} + y^{28} = z^{28}$.
7. Find all Pythagorean triangles with perimeter 1716.

What Is a Calculation?

At the second International Congress of Mathematicians (Paris, 1900), David Hilbert (1862–1943) presented a list of 23 problems, which he hoped would occupy mathematicians in the 20th century. We shall only talk about three of these problems here, as they concern the foundations of mathematics.

1. Prove or disprove the Continuum Hypothesis (Chapter 12).
2. Show that arithmetic, described as an axiomatic system, is consistent. that is, that it does not admit a proof that $0 = 1$. (As Paul Erdős would say, if such a proof were ever to be discovered, the universe would vanish.)
3. Find an effective method or ‘algorithm’, as it is now called, for deciding whether a given polynomial Diophantine equation (with integer coefficients) is solvable (in integers). (For the origin of the word ‘algorithm’, see Part I, Chapter 22.)

This last problem is actually Hilbert’s Problem 10. Today we know that these three problems cannot be solved in the way Hilbert had intended. Kurt Gödel showed in 1938 that the Continuum Hypothesis cannot be proved and Paul Cohen showed in 1964 that it cannot be disproved either! The existence of mathematical statements that can be neither proved nor disproved had already been established by Gödel in 1931. It followed from his argument that the consistency of any formal system of arithmetic cannot be proved, unless we allow a method of proof which is more powerful than

the method of proof in the given system. We shall return to Gödel's result in a later chapter. After much preliminary work by Martin Davis, Hilary Putnam and Julia Robinson, Problem 10 was ultimately disposed of by Yuri Matijasevič, as we shall see presently. But before tackling this problem, one had to determine what Hilbert meant by an 'effective method'.

By a *numerical function* we mean any function $f : \mathbf{N}^n \rightarrow \mathbf{N}$, where \mathbf{N} is the set of natural numbers (including 0) and $n \geq 0$. Among these are the *identity* function $fx = x$, and the successor function $fx = Sx$. These basic functions can surely be calculated. Moreover, the set of calculable functions is evidently closed under the following schemes, where we have written \bar{x} for the string $x_1x_2 \dots x_m$, \bar{z} for $z_1 \dots z_k$, etc:

1. substituting one calculable function $g\bar{y}$ for u in another calculable function $f\bar{x}u\bar{z}$ to obtain a function $h\bar{x}\bar{y}\bar{z} = f\bar{x}g\bar{y}\bar{z}$;
2. interchanging two variables: if $f\bar{u}xy\bar{v}$ can be calculated, so can $g\bar{u}xy\bar{v} = f\bar{u}yx\bar{v}$;
3. contracting two variables: if $f\bar{u}xy\bar{v}$ can be calculated, so can $g\bar{u}x\bar{v} = f\bar{u}xx\bar{v}$;
4. introducing superfluous variables: if $f\bar{u}\bar{v}$ can be calculated, so can $g\bar{u}x\bar{v} = f\bar{u}\bar{v}$;
5. the *recursion scheme*: if $g\bar{x}$ and $h\bar{x}yz$ are calculable, so is $f\bar{x}y$ defined 'recursively' by the equations

$$f\bar{x}0 = g\bar{x}, \quad f\bar{x}Sy = h\bar{x}yf\bar{x}y.$$

With the help of (5) we can calculate $x + y$, $x \times y$, x^y and many other numerical functions.

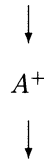
The functions which can be calculated so far, namely from the basic functions with the help of schemes (1) to (5), are called *primitive recursive functions*. They were introduced by Gödel in the proof of his famous Incompleteness Theorem, which led to the disposal of Hilbert's second problem and to which we shall return in a later chapter.

How are these primitive recursive functions to be calculated? Before the invention of modern computers and pocket calculators, we could make calculations with pencil and paper; but to explain what this means in precise terms is not easy. The first people to give rigorous answers to the question of what a calculation is were Alan Turing and Emil L. Post, independently in the same year, 1936. The work of Post is not as well-known as that of Turing, who invented a theoretical machine, the *Turing machine*, which may be seen as the ancestor of all modern computers. However, building a Turing machine is not an easy way to compute a given primitive recursive function.

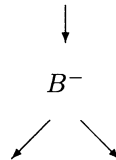
A simpler answer would have been available if people had thought about the origin of the word ‘calculation’, which is derived from the Latin word ‘calculus’ meaning ‘pebble’. Indeed, the ancients performed calculations by moving pebbles from one groove of a table, called an ‘abacus’, to another. We shall discuss how an abacus may be programmed to calculate any primitive recursive function.

For theoretical purposes we shall assume that an *abacus* consists of a potentially infinite number of *locations* and that we have an inexhaustable supply of *pebbles*. A *program* is made up of two basic instructions:

- put a pebble at location A , then go to ...

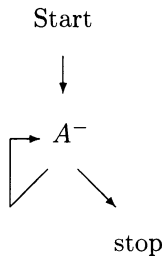


- take a pebble from location B , if B is not empty, then go to ..., else go to ...

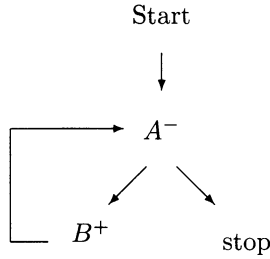


(the right arrow referring to the ‘else’ case).

A program is easily illustrated by a *flow diagram*. For example, the following program

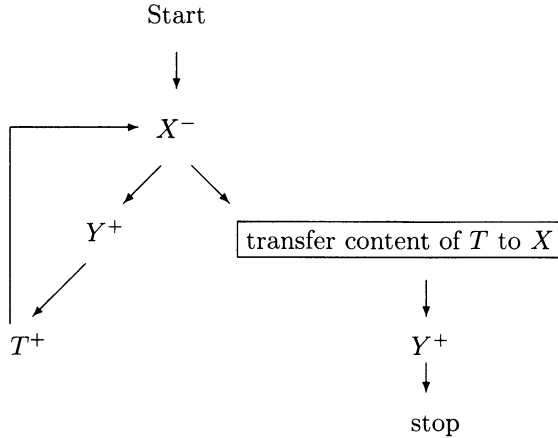


is the program ‘empty A ’ and the next program



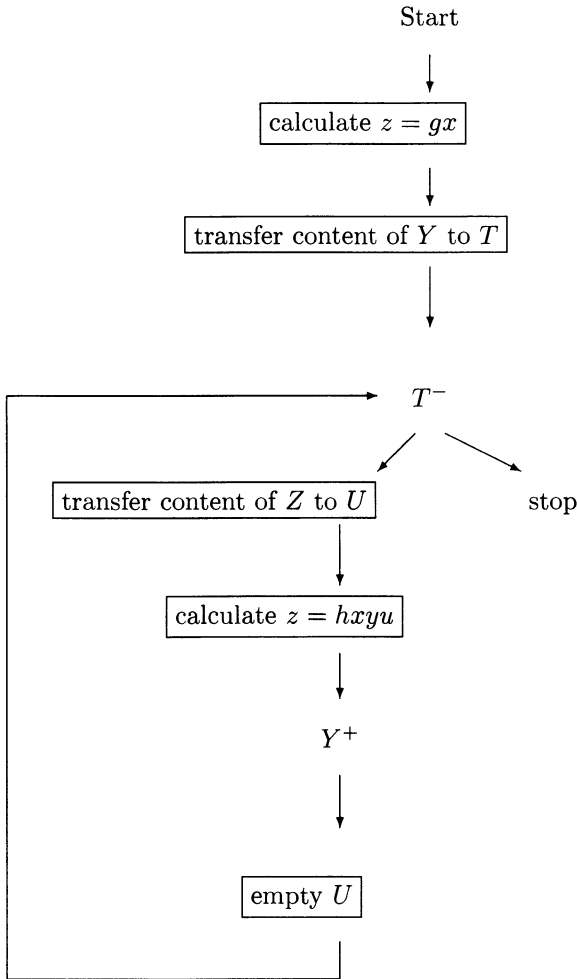
is the program ‘transfer the content of A to B ’.

To *calculate* a numerical function $z = fxy$ we begin with x pebbles at location X , y pebbles at location Y and no pebbles at any other location. We expect the calculation to stop with x pebbles at X , y pebbles at Y , fxy pebbles at Z and no pebbles at any other location. Similarly we define what it means to calculate $z = fx_1, \dots, x_n$. The following flow diagram illustrates the programs for calculating the successor function $y = Sx$:



Note the *subroutine* ‘transfer content of T to X ’. T is a temporary storage location, which is empty at the beginning and at the end of the calculation.

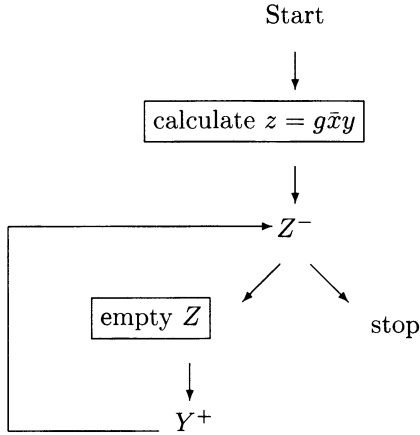
To convince the reader that in fact all primitive recursive functions can be calculated on an abacus, we present a program for calculating a function fxy obtained by the recursion scheme $fx0 = gx$, $fxSy = hxyfxy$ from functions gx and $hxyz$, which are already known to be calculable:



Primitive recursive functions are not the only numerical functions which can be calculated on an abacus. There is one more scheme, the *minimization scheme*:

6. given a calculable function $g\bar{x}y$, $\bar{x} = x_1x_2 \dots x_n$, we can calculate $f\bar{x} = \text{the smallest } y \text{ such that } g\bar{x}y = 0$.

Here is the flow diagram:



Of course, $f\bar{x}$ is only well-defined if $\forall x\exists y(g\bar{x}y = 0)$, otherwise f is not a function but a *partial function*. The class of (partial) functions obtained by adding the minimization scheme to the schemes discussed earlier is the class of *(partial) recursive functions*. These are in fact the only numerical (partial) functions which can be calculated on an abacus or, for that matter, on a Turing machine or on a modern computer.

Exercises

Construct programs for calculating the following numerical functions on an abacus:

1. $z = x + y, x \cdot y, x^y$;
2. $z = y!$;
3. $y = \lfloor x/2 \rfloor$, where $\lfloor \alpha \rfloor$ is the greatest integer in α ;
4. $y = \lfloor \sqrt{x} \rfloor$.

Recursive and Recursively Enumerable Sets

A set A of natural numbers is said to be *recursive* if there is a recursive (or calculable) function $f : \mathbf{N} \rightarrow \mathbf{N}$ such that $A = \{x \in \mathbf{N} | f(x) = 0\}$. In other words, A is recursive provided, for any natural number x , we can determine whether $x \in A$ by performing a calculation f on x . (Actually, it suffices here to take for f a primitive recursive function.) For example, the set of primes is recursive, since we can determine whether a given number is prime by dividing it by all natural numbers less than or equal to its square root. With a bit of work, one can find a recursive function f to do the job.

A set of natural numbers is *recursively enumerable* if there is a recursive function $g : \mathbf{N} \rightarrow \mathbf{N}$ such that $A = \{g(0), g(1), g(2), \dots\}$, where the $g(n)$ are not necessarily arranged in order of magnitude and may repeat. (Again, it actually suffices to take g to be primitive recursive.) In other words, A is recursively enumerable if there is a calculation which will generate all and only the elements of A in some order, possibly with repetitions. Although the empty set does not conform to the definition above, it will be convenient to consider it to be recursively enumerable.

An important observation linking the above two concepts is the following due to Kleene:

Proposition 20.1. *Let A be a set of natural numbers and A^c its complement in \mathbf{N} . Then A is recursive if and only if both A and A^c are recursively enumerable.*

Proof sketched: Suppose A is recursive. Then there is a calculable function f such that $x \in A$ if and only if $f(x) = 0$. We may assume that A is not empty, for otherwise the result holds trivially. Let a be the smallest element

of A . Define $g(x) = x$ if $f(x) = 0$ and $g(x) = a$ if $f(x) \neq 0$. Then g is surely calculable and A is its range. Thus A is recursively enumerable. One shows similarly that A^c is recursively enumerable.

Conversely, suppose both A and A^c are recursively enumerable, both nonempty. Then A is the range of a calculable function g and A^c is the range of a calculable function h . Now every natural number is in the range of g or the range of h and we can determine whether x is an element of A by calculating

$$g(0), h(0), g(1), h(1), g(2), \dots$$

We put $f(x) = 0$ if, for some y , $g(y) = x$ and $f(x) = 1$ if, for some y , $h(y) = x$. Then f is calculable and $A = \{x \in N \mid f(x) = 0\}$, hence A is a recursive set.

The above considerations can be extended from natural numbers to finite strings of symbols in a formal language, e.g., the language of mathematics. We assume that the set of symbols is finite. For example, if mathematics is based on set theory, we know that the following symbols suffice:

$$x, ' (,), \wedge, \Rightarrow, \vee, \neg, \forall, \exists, =, \in, 0, S.$$

Assuming that there are s symbols in the alphabet (e.g., $s = 14$), we can think of a string of symbols as a natural number expressed in the base s . If you prefer, you may say that the number *encodes* the string.

A mathematical *formula* is a finite string of symbols, combined according to some simple rules, and so is a *proof* in mathematics: a finite list of formulas, each of which is an axiom or else derived from earlier formulas in the list by a rule of inference. There is a finite procedure for 'calculating' whether a formula in a given list meets one of these two requirements, hence whether the given list is a proof. It follows that the set of proofs in mathematics, viewed as the set of natural numbers which encode these proofs, is a recursive set. The last formula in a proof is called a *theorem*.

We can program a computer to look at each natural number in turn, convert it to base s , and determine whether it is (encodes) a formula or a proof. In the second case, we can also tell the computer to print out the number which encodes the last formula of the proof. The computer will then generate a list which consists of all and only those base s natural numbers which encode theorems of mathematics. It follows that the set of theorems in mathematics is recursively enumerable.

On the other hand, it was proved by Alonzo Church that there is no mechanical procedure or algorithm for determining, in general, whether a given formula is a theorem of mathematics. In other words, the set of theorems is not recursive.

If we combine this result of Church with the above proposition of Kleene, we may deduce that the complement of the set of theorems is not recursively enumerable. Moreover, we may deduce the following:

Theorem 20.2. (Incompleteness Theorem)

There are mathematical formulas p such that neither p nor $\neg p$ are provable.

Proof: Let T be the set of theorems and C the set of contradictions, that is, formulas p such that $\neg p$ is a theorem. The computer can print out the members of C as easily as those of T , hence both sets are recursively enumerable. Suppose that, for every formula p , either p is provable or $\neg p$ is provable. Then p is a nontheorem if and only if $\neg p$ is provable, hence C is the complement of T . Thus both T and its complement are recursively enumerable, which contradicts the result of Church.

We should note that the incompleteness theorem was first proved by Gödel (1906–1978) and that Church used Gödel's result to prove his own. In a subsequent chapter, we shall present Gödel's original proof, which does not depend on Church's result.

Hilbert's Tenth Problem

Hilbert's tenth problem asked for an algorithm to determine whether any given polynomial Diophantine equation has a solution in integers. After important preliminary work by Martin Davis, Hilary Putnam (the philosopher) and Julia Robinson, Yuri Matiyasevič showed that no such algorithm exists. The proof is long and we shall give only a few of the highlights here. For a complete treatment, the reader may wish to consult Davis [1973].

First let us reduce the problem of solving a Diophantine equation in *integers* to one of solving it in *positive integers*, using the fact that every positive integer x can be written in the form

$$x = x_0^2 + x_1^2 + x_2^2 + x_3^2 + 1,$$

where the x_i are integers, in view of Lagrange's Theorem (Chapter 9). For example, if we want to know whether

$$x^{17} + y^{17} - z^{17} = 0$$

has a solution in positive integers, we may test whether the following equation has a solution in integers:

$$(x_0^2 + x_1^2 + x_2^2 + x_3^2 + 1)^{17} + (y_0^2 + \cdots + 1)^{17} - (z_0^2 + \cdots + 1)^{17} = 0.$$

A set A of positive integers is said to be *Diophantine* if there is a polynomial $p(t, x_1, \dots, x_n)$ with integer coefficients such that $t \in A$ if and only if there are positive integers x_1, x_2, \dots, x_n such that $p(t, x_1, \dots, x_n) = 0$. We shall write $A = A_p$ to express the relationship between the set A and

the polynomial p . For example, the set of *composite* positive integers is Diophantine of the form A_p , where

$$p(t, x_1, x_2) \equiv t - (x_1 + 1)(x_2 + 1) \equiv t - x_1x_2 - x_1 - x_2 - 1.$$

Lemma 21.1. *Every Diophantine set is recursively enumerable.*

Proof: Given the polynomial $p(t, x_1, \dots, x_n)$ with integer coefficients and any positive integer m , let S_m be the set of all $(n+1)$ -tuples of positive integers, each less than or equal to m . Then S_m is clearly finite. We can enumerate the elements of the Diophantine set A_p by looking at each of S_1, S_2, \dots in turn and checking whether it contains a solution of the equation $p(t, x_1, x_2, \dots, x_n) = 0$. Whenever we do find a solution, we list its first member t .

Lemma 21.2. *Suppose there is an algorithm for deciding whether, for any given t , the polynomial equation $p(t, x_1, \dots, x_n) = 0$ has a solution in positive integers. Then the Diophantine set A_p is recursive.*

Proof: This is so because we can perform a calculation on t to see whether $t \in A_p$.

The next lemma is due to Julia Robinson. It has to do with the notion of *exponential growth*. For example, consider the *Fibonacci sequence* F_m :

$$1, 1, 2, 3, 5, 8, 13, \dots$$

in which each term, except $F_1 = 1$ and $F_2 = 1$, is the sum of the preceding two terms: $F_{m+2} = F_{m+1} + F_m$. This had been studied by Leonardo of Pisa (1180–1250), also known as ‘Fibonacci’, in connection with the growth of rabbit populations. As we saw in Part I, Chapter 23,

$$F_m = \frac{(\frac{1}{2}(1 + \sqrt{5}))^m - (\frac{1}{2}(1 - \sqrt{5}))^m}{\sqrt{5}}.$$

Since $(\frac{1}{2}(1 - \sqrt{5}))^m / \sqrt{5}$ is small, F_m is in fact equal to the integer nearest to $(\frac{1}{2}(1 + \sqrt{5}))^m$. This explains why rabbit populations grow exponentially.

Lemma 21.3. (Julia Robinson)

A sufficient condition for every recursively enumerable set to be Diophantine is that there is a polynomial equation with integer coefficients

$$p(u, v, x_2, \dots, x_n) = 0$$

such that, in its positive integer solutions, v grows exponentially relative to u .

The proof of Lemma 21.3 is too long to be included here.

In 1970, Matiyasevič found a polynomial Diophantine equation whose $(m - 1)$ th positive integer solution, for $m \geq 2$ has the form

$$(m, F_{2m}, x_2, \dots, x_n).$$

By Lemma 21.3, he inferred

Theorem 21.4. (Matiyasevič)

Every recursively enumerable set of positive integers is Diophantine.

It follows by Lemma 21.1 that a set of positive integers is Diophantine if and only if it is recursively enumerable. We saw in Chapter 20 that not every recursively enumerable set of positive integers is recursive, e.g., the set of positive integers which encode the theorems of mathematics. Therefore, in view of Lemma 21.2, there is no algorithm for testing whether any given polynomial equation has a solution in positive integers. We conclude that *Hilbert's tenth problem is unsolvable*.

Matiyasevič's result had curious repercussions on the existence of *prime representing polynomials*. Consider the polynomial

$$f(x) = x^2 - x + 41.$$

For $x = 1, 2, 3, \dots, 40$, $f(x)$ is a prime number. While this might convince a physicist that, for any positive integer x , $f(x)$ is always prime, $x = 41$ is a counterexample.

On the other hand, consider the polynomial

$$g(x) = -x^2 + 3.$$

The set of all $g(x) \geq 0$ such that x is a positive is a *subset* of the set of prime numbers, albeit a very small subset.

These examples suggest the following questions:

1. Is there a polynomial with integer coefficients all of whose values for positive integer arguments are primes?
2. Is there a polynomial with integer coefficients such that the set of its *positive* values for positive integer arguments is just the set of primes?

The answer to the first question is 'no', as the reader will be invited to verify in the Exercises. Surprisingly, the answer to the second question is 'yes'. This uses nothing about the set of prime numbers except that it is recursively enumerable, hence Diophantine by Theorem 21.4.

Theorem 21.5. (Putnam)

For any Diophantine set A_p there is a polynomial

$$q(t, x_1, x_2, \dots, x_n)$$

with integer coefficients such that A_p is the set of all positive values of $q(t, x_1, \dots, x_n)$ for positive integers t, x_1, \dots, x_n .

Proof: We recall that $t \in A_p$ if and only if $p(t, x_1, \dots, x_n) = 0$. Let

$$q(t, x_1, \dots, x_n) = t(1 - (p(t, x_1, \dots, x_n))^2).$$

We shall illustrate the argument by taking $n = 1$.

(a) Suppose t and x are positive integers and $q(t, x) > 0$. Then

$$t(1 - (p(t, x))^2) > 0$$

and hence $1 > (p(t, x))^2$, so $p(t, x) = 0$, and hence $t \in A_p$.

(b) Suppose $t \in A_p$. Then there is a positive integer x such that $p(t, x) = 0$, hence $q(t, x) = t$.

It follows from Putnam's Theorem that there is a polynomial

$$q(t, x_1, \dots, x_n)$$

with integer coefficients, such that the set of its positive values, for positive integer assignments of the $n + 1$ variables, is exactly the set of prime numbers. Such a polynomial can be found in Browder [1976], p. 331.

Exercises

1. Explain why the set of positive integers which are not powers of 2 is Diophantine.
2. Describe an algorithm for solving Diophantine equations in one variable:

$$a_0x^n + a_1x^{n-1} + \dots + a_n = 0,$$

where the a_i are given integers.

3. Let $f(x)$ be a polynomial with integer coefficients. Show that there is a positive integer x such that $f(x)$ is not prime.
4. Let $f(x_0, x_1, \dots, x_n)$ be a polynomial with integer coefficients. Show that there are positive integers x_0, \dots, x_n such that $f(x_0, \dots, x_n)$ is not prime.
5. Prove that the set $\{2, 3, 4\}$ is Diophantine and find a polynomial $q(t, x_1, \dots, x_n)$ with integer coefficients whose positive values, for positive integer arguments, are just the members of this set.
6. Find a polynomial with integer coefficients whose positive values, for positive integer arguments, are all and only composite numbers.

Lambda Calculus

The Lambda Calculus of Alonzo Church represents an attempt to understand mathematical entities as *functions*. Usually, people think of a function $f : A \rightarrow B$ as having a domain A and a codomain B . But, in the *untyped* version of the lambda calculus, one makes the implicit assumption that $A = B$ is some kind of universal set and that f is defined everywhere; it is even possible to apply f to itself.

We write $f'a$ for the value of f at a and we read this as f of a . For example, if f is the squaring function, we write $f'a$ for a^2 . Similarly, if $\phi(x)$ is the expression $x^3 + x + 1$, we can introduce a function g such that $g'x = \phi(x) \equiv x^3 + x + 1$. It is customary to write $g = \lambda_x(x^3 + x + 1)$, where λ_x is the *abstraction operator*. It is sometimes important to distinguish between the expression $\phi(x)$ and the function $\lambda_x\phi(x)$ which sends x to $\phi(x)$. In particular, $\lambda_x\phi(x)'2 = 2^3 + 2 + 1 = 11$ and $\lambda_x\phi(x)'y = \phi(y) \equiv y^3 + y + 1$. In 1937, Church and Turing showed independently that every calculable numerical function can be expressed in terms of the untyped lambda calculus. In particular, the natural numbers and the usual arithmetical operations on natural numbers can be so expressed, as we shall see.

Conversely, every numerical function definable in terms of the untyped lambda calculus is calculable, thus

$$\text{recursive} = \text{calculable} = \lambda\text{-definable}.$$

According to the *Church-Turing Thesis*, any of these three equivalent concepts captures the intuitive notion of what it means to be ‘computable’.

We shall now give a rigorous presentation of the *untyped lambda calculus*. We assume, to begin with, that there is a supply of countably many

variables x_1, x_2, \dots . Rather than referring to specific variables, such as x_{17} , we shall use letters x, y, z, f, g, h, \dots as arbitrary variables. We now define *terms* of the λ -calculus:

1. each variable is a term;
2. if ϕ and ψ are terms, so is $(\phi'\psi)$;
3. if ϕ is a term and x is a variable, then $(\lambda_x\phi)$ is a term;
4. all terms are built up according to the above rules.

We have inserted parentheses to avoid ambiguity. However, when there is no danger of ambiguity we may omit them.

We define a *free* occurrence of a variable in a term as follows:

1. a variable x occurs freely in the term x ;
2. if x occurs freely in ϕ or ψ , then it occurs freely in $\phi'\psi$;
3. if y is a variable distinct from x and y occurs freely in ϕ , then it also occurs freely in $\lambda_x\phi$; however, x does *not* occur freely in $\lambda_x\phi$.

If the variable x occurs freely in the term ϕ , then it is said to be *bound* to λ_x in $\lambda_x\phi$. We often write the term ϕ as $\phi(x)$ to indicate the possible free occurrence of x in ϕ . This has the advantage that we can write $\phi(\alpha)$ for the result of substituting α for x in ϕ . But we are only *allowed to substitute* α for x in $\phi(x)$ if no variable occurring freely in α becomes bound in $\phi(\alpha)$.

In addition to the usual properties of equality — reflexivity, symmetry, transitivity and the rule allowing substitution of equals — we postulate the following three rules:

- R1.** $(\lambda_x\phi(x))'\alpha = \phi(\alpha)$, if we are allowed to substitute α for x in $\phi(x)$;
- R2.** $\lambda_x(\phi'x) = \phi$, if x is not free in ϕ (possibly because x does not occur in ϕ at all);
- R3.** $\lambda_x\phi(x) = \lambda_y\phi(y)$, if we are allowed to substitute y for x in $\phi(x)$.

Here are some examples to illustrate these ideas, which the reader may skip if she has already understood them.

1. $(\lambda_x(y'x))'(x'z)$ is a term with the first occurrence of x not free, it being bound to λ_x , and the second occurrence of x free.

2. Suppose $\phi(x)$ and $\psi(x)$ are terms not containing λ_y and suppose we know that $\phi(x) = \psi(x)$. Then $\lambda_x\phi(x)'y = \lambda_x\psi(x)'y$, by the properties of equality, hence $\phi(y) = \psi(y)$ by R1.
3. $\lambda_f(g'f) = g$ by R2 (assuming $g \neq f$) and $(\lambda_gg)'f = f$ by R1.
4. What can go wrong if we disregard the 'if' clauses of the rules? From 3 we see that

$$(\lambda_g(\lambda_f(g'f)))'f = (\lambda_gg)'f = f.$$

But disregarding the 'if' clause in R1 would yield

$$(\lambda_g(\lambda_f(g'f)))'f = \lambda_f(f'f)$$

(taking $x \equiv g$ and $\phi(x) \equiv \lambda_f(g'f)$). But f and $\lambda_f(f'f)$ are not the same, as is seen by applying both to a , say. $f'a$ depends on f , but $(\lambda_f(f'f))'a = a'a$ does not; f in the last term is a *dummy* variable.

We shall now show how to do arithmetic in the lambda calculus. First note that any two functions f and g can be *composed* to form $f \circ g$, where $(f \circ g)'x = f'(g'x)$, hence $f \circ g = \lambda_x(f'(g'x))$ by R2. In particular, $f \circ f$ is usually written f^2 , the *iterate* of f , so $f^2 = \lambda_x(f'(f'x))$. Church had the idea that the number 2 should be defined as the *process of iteration*, as the function which assigns f^2 to f , that is

$$2 = \lambda_f(f^2) = \lambda_f(\lambda_x(f'(f'x))).$$

Since it is customary to think of f^0 as the identity function and of f^1 as f , we may also define

$$0 = \lambda_f(\lambda_x x), \quad 1 = \lambda_f(\lambda_x(f'x)) = \lambda_f f,$$

hence $0 = \lambda_f 1$ since $\lambda_f f = \lambda_x x$ by R3. The reader is invited to define the number 3 in the lambda calculus. In general, $n = \lambda_f(f^n)$, where $f^n = f \circ f \circ \dots \circ f$ is the n th iterate of f . More precisely, f^n is defined inductively by $f^{\Sigma'n} = f \circ f^n$, where $\Sigma'n$ is the successor of n . Thus $n'f = f^n$.

Substituting m for f in this equation, we have $n'm = m^n$. But this tells us how to define *exponentiation* in the lambda calculus. In particular, $m^0 = 0'm = 1$. Since $(m \times n)'f = f^{m \times n} = (f^n)^m = m'(n'f) = (m \circ n)'f$, we may define *multiplication* by $m \times n = m \circ n = \lambda_x(m'(n'x))$. For *addition* we have

$$(m + n)'f = f^{m+n} = f^m \circ f^n = (m'f) \circ (n'f),$$

so that $m + n = \lambda_f((m'f) \circ (n'f))$.

We can now prove some of the basic laws of arithmetic. For example, since composition of functions is associative, multiplication of numbers is associative. On the other hand, composition of functions is not commutative, yet multiplication of numbers is; this must be proved by mathematical induction.

Of considerable interest in computer science is the following:

Theorem 22.1. (Fixpoint Theorem)

For any term ϕ , there is a term α such that $\phi'\alpha = \alpha$ can be proved.

This implies, in particular, that it is not possible to incorporate a term into the lambda calculus which corresponds to the negation symbol in logic.

Proof: Let $\beta = \lambda_x(\phi'(x'x))$ and $\alpha = \beta'\beta$. Then

$$\alpha = \beta'\beta = (\lambda_x(\phi'(x'x)))'\beta = \phi'(\beta'\beta) = \phi'\alpha.$$

Surprisingly, it is possible to get rid of the λ -abstraction (and all bound variables) in the lambda calculus. This discovery goes back to Schönfinkel (1924). In fact, if we write

$$\begin{aligned} I &= \lambda_x x && \text{(identity),} \\ K &= \lambda_x \lambda_y x && \text{(constancy operator),} \\ S &= \lambda_x \lambda_y \lambda_z ((x'z)'(y'z)) && \text{(Schönfinkel operator),} \end{aligned}$$

we have

$$\begin{aligned} I'x &= x && \text{(thus } I = 1), \\ (K'x)'y &= x && \text{(thus } K'x \text{ is the function} \\ &&& \text{with constant value } x), \\ ((S'f)'g)'x &= (f'x)'(g'x). \end{aligned}$$

We define *combinators* as follows:

1. all variables are combinators;
2. I , K , and S are combinators;
3. if ϕ and ψ are combinators, so is $\phi'\psi$;
4. all combinators are obtained from the above rules.

Theorem 22.2. (Schönfinkel)

Every term of the lambda calculus is provably equal to a combinator.

Proof: In view of rules (1) to (3), this will follow if we show that, if ϕ is equal to a combinator, so is $\lambda_x \phi$. We prove this by induction on the length of ϕ .

Since $\lambda_x x = I$, $\lambda_x y = K'y$, $\lambda_x I = K'I$, $\lambda_x K = K'K$ and $\lambda_x S = K'S$, we know that the induction hypothesis holds for all combinators ϕ of length 1.

Suppose now that ϕ has length greater than 1 and that every combinator ψ shorter than ϕ is such that $\lambda_x \psi$ is equal to a combinator. We claim that $\lambda_x \phi$ is equal to a combinator. Indeed, $\phi = \psi' \chi$, where ψ and χ are shorter than ϕ , hence

$$\lambda_x \phi = \lambda_x (\psi' \chi) = (S'(\lambda_x \psi))'(\lambda_x \chi).$$

This is a combinator because $\lambda_x \psi$ and $\lambda_x \chi$ are combinators by our inductive assumption.

Using the methods of this proof, we can express 2 as a combinator:

$$\begin{aligned} 2 &= \lambda_f \lambda_x (f'(f'x)) \\ &= \lambda_f ((S' \lambda_x f)' \lambda_x (f'x)) \\ &= \lambda_f ((S' \lambda_x f)' f) && \text{by R2} \\ &= (S'(\lambda_f (S' \lambda_x f)))' \lambda_f f \\ &= (S'(\lambda_f (S'(K'f))))' I \\ &= (S'((S' \lambda_f S)' \lambda_f (K'f)))' I \\ &= (S'((S'(K'S))'K))' I && \text{by R2.} \end{aligned}$$

Exercises

1. Assuming that m, p and q are natural numbers expressed in the lambda calculus, show that $(m^p)^q = m^{(q \times p)}$.
2. Prove that $0^{\Sigma' n} = 0$.
3. Prove that $I = (S'K)'K$.
4. Express m^n , $m \times n$, and $m + n$ in terms of I, K, S, m and n .

23

Logic from Aristotle to Russell

Logic was not always regarded as a branch of mathematics, certainly not by Aristotle (384–322 BC), who was the first to write about logic in the West. Among the principles which he recognized are the following:

$$\begin{array}{ll}\neg\neg p \iff p & \text{(double negation),}\\ p \vee \neg p & \text{(excluded third),}\\ (p \Rightarrow q) \iff (\neg q \Rightarrow \neg p) & \text{(contraposition).}\end{array}$$

He also looked at modal logic and showed how possibility can be defined in terms of necessity.

Aristotle's major concern was with a type of argument called the 'syllogism', which predominated in logical thinking for the next two thousand years. It dealt with four types of basic statements:

SaP	meaning <i>all S are P</i> ,
SeP	meaning <i>no S are P</i> ,
SiP	meaning <i>some S are P</i> ,
SoP	meaning <i>some S are not P</i> .

He realized that PeS is equivalent to SeP and that PiS is equivalent to SiP and he adhered to a convention that SaP implies SiP. (Today we use

words differently: we assert that *all* unicorns have horns, but deny that *some* unicorns have horns. Evidently Aristotle did not believe in the empty set.)

A *syllogism* is an argument which infers one such basic statement from two others. Here are the first four ‘figures’ of the syllogism:

MaP	MeP	MaP	MeP
SaM	SaM	SiM	SiM
<hr/> SaP	<hr/> SeP	<hr/> SiP	<hr/> SoP

William of Shyreswood (1250 AD) gave these syllogisms the names

barbara, celarent, darii, ferio,

— to make them easier to remember. There were more such figures, which we shall not discuss here. Here is a typical argument illustrating the ‘ferio’:

no minister is prudent
some socialists are ministers
<hr/>
some socialists are not prudent

The Stoics (200 BC), Philo of Megara in particular, essentially introduced truth tables into logic, thus anticipating Ludwig Wittgenstein (1889–1951). They discussed the problem of whether ‘*p* or *q*’ is true when both *p* and *q* are true and whether ‘if *p* then *q*’ is always true when *p* is false. They arrived at the modern conventions, expressed in the following *truth tables*:

$p \vee q$	$p \Rightarrow q$
<hr/>	<hr/>
T T T	T T T
T T F	T F F
F T T	F T T
F F F	F T F

The Stoics were committed to the view that there are only two truth values (T and F). In particular, they believed that a statement like ‘there will be a battle tomorrow’ is either true or false, although Aristotle seems to have had some second thoughts about this. This belief was associated in

their minds with the view that the future is determined by fate (Kneale p. 48).

Today we tend to dismiss the detailed elaboration of Aristotelian logic by the medieval Scholastics and recognize as a major advance only the ideas of Gottfried W. Leibniz (1646–1716). Leibniz conceived of a universal symbolic language which would be adequate not only for mathematics, but for all of science. Unfortunately, he did not get around to publishing the details of his proposal, perhaps because he was preoccupied by his controversy with Newton concerning the invention of ‘the’ calculus and by his successful diplomatic efforts to put George I on the throne of England.

It was only in 1847 that full-blown *symbolic logic* finally saw the light of day. This was the year in which both George Boole (1815–1864) and Augustus DeMorgan (1806–1871) published their first works in logic. The former saw propositional logic as a branch of algebra, distinguished from the usual algebra of ‘quantities’ by the ‘idempotent’ law: $p \times p = p$. De Morgan is remembered for his laws expressing the duality between disjunction (‘or’) and conjunction (‘and’):

$$\neg(p \wedge q) = \neg p \vee \neg q, \quad \neg(p \vee q) = \neg p \wedge \neg q.$$

The next major step was taken by Gottlob Frege (1848–1925), who was the first to have a modern view of universal and existential quantifiers, although without using the modern notation: \forall and \exists . This was surprisingly recent, considering that every student nowadays is familiar with these notions. Frege also attempted to express all of mathematics in terms of logical symbols, in fact, reducing mathematics to logic, thus espousing a philosophical position called ‘Logicism’.

Crucial to Frege’s project is the assumption that, corresponding to any property expressed by a predicate of his language, there is a uniquely determined set whose elements are just the entities with that property. Frege expressed this assumption by the following *comprehension scheme*:

$$\exists y \forall x (x \in y \iff P(x)),$$

where $P(x)$ is any formula, possibly containing the free variable x . (The double arrow means ‘if and only if’ or ‘just in case’.) This scheme is accompanied by the *axiom of extensionality* to ensure the uniqueness of the set y whose existence has been asserted above:

$$\forall y \forall z (\forall x (x \in y \iff x \in z) \implies y = z).$$

The axiom of extensionality implies, in particular, that there can be at most one entity y with no elements, the *empty set*. This axiom implicitly contains the assumption that all entities discussed by Frege’s language are sets.

Frege had just written a book propounding these views when he received a letter from Bertrand Russell (1872–1970), pointing out that there was a

serious problem when $P(x)$ was the formula $\neg(x \in x)$. Indeed, if y was such that $\forall(x \in y \iff \neg(x \in x))$, one would obtain as a special case

$$y \in y \iff \neg(y \in y),$$

which is a contradiction.

This argument is known as ‘Russell’s paradox’. One way to avoid the contradiction is to forbid expressions such as $x \in x$. Russell and Whitehead propose a *theory of types*, according to which each symbol denoting an entity should have attached to it a certain natural number, its *type*, and the formula $a \in b$ is permitted only if the type of b is one higher than the type of a . This theory was developed in excruciating detail in the three-volume *Principia Mathematica*, an unnecessarily complicated treatise that is more talked about than studied.

Although up-to-date versions of type theory are now available, most mathematicians prefer other methods for avoiding Russell’s and similar paradoxes. On the whole, if mathematicians worry about such problems at all, they subscribe to the set theory of Gödel and Bernays, which distinguishes between *sets* and *classes*: only the former can be elements. Unfortunately, one has to add a number of axioms to specify which classes are sets. Logicians, on the other hand, prefer the set theory of Ernst Zermelo (1871–1953) and Abraham Fraenkel (1891–1965), who modify the comprehension as follows:

$$\forall z \exists y \forall x (x \in y \iff (x \in z \wedge P(x))).$$

They too had to introduce additional axioms, which spoiled the simplicity of Frege’s project.

The fact that Frege’s simple and natural comprehension scheme led to contradictions startled many mathematicians, and some became sceptical about all but the most basic procedures. Other mathematicians were already sceptical. We shall discuss L. Kronecker (1823–1891) and H. Poincaré (1854–1912) in this chapter, leaving L. E. J. Brouwer to the next chapter.

It was Kronecker who said ‘God made the whole numbers, all the rest is the work of man’. He was suspicious of Cantor’s infinite cardinals; but most of all he rejected *nonconstructive* arguments such as the following proof that there exist irrational numbers α and β such that α^β is rational.

Consider $\sqrt{2}^{\sqrt{2}}$. If it is rational, we are done, since we know that $\sqrt{2}$ is irrational. Suppose $\sqrt{2}^{\sqrt{2}}$ is irrational. Call it α and take $\beta = \sqrt{2}$. Then

$$\alpha^\beta = \sqrt{2}^{(\sqrt{2} \times \sqrt{2})} = \sqrt{2}^2 = 2,$$

which is surely rational.

This proof depends on the Stoic idea, also endorsed by Aristotle, that every proposition, here the proposition that $\sqrt{2}^{\sqrt{2}}$ is rational, is either true or false. What Kronecker would have objected to is that, at the end of the proof, we don't know whether $\alpha = \sqrt{2}^{\sqrt{2}}$ or $\alpha = \sqrt{2}$.

It is known, using some deep mathematics, that $\sqrt{2}^{\sqrt{2}}$ is irrational, but that is beside the point here. It is also beside the point that the stated theorem has an easy constructive proof: take $\alpha = \sqrt{2}$ and $\beta = 2\log_2 3$. It is possible to exhibit mathematical theorems for which no constructive proof exists, for instance the theorem which asserts that, if the axiom of choice is true, then there exist nonprincipal ultrafilters on the set of natural numbers. (The reader unfamiliar with ultrafilters may ignore this example.)

Poincaré, on the other hand, objected to *impredicative* definitions or constructions, essentially those which define or construct an entity in terms of entities of a higher type. For example, consider the usual proof that every nonempty set of real numbers that is bounded above has a least upper bound. According to Dedekind, a *real number* is a set α of rational numbers such that

1. both α and its complement α^c are nonempty;
2. every element of α is less than every element of α^c ;
3. α has no greatest element.

Now let m be a nonempty set of real numbers which is bounded above. Put

$$\alpha \equiv \{x \in \mathbf{Q} \mid \exists y \in m, x \in y\}.$$

It is not hard to show that α satisfies 1 to 3 and that it is the least upper bound of m . However, being a *set* of reals, m has a higher type than α . The construction of α is thus impredicative in Poincaré's sense.

Most mathematicians feel that Poincaré was too sceptical here. If we disallowed constructions such as that of the least upper bound of m above, most of analysis would have to be abandoned.

Exercises

1. Suppose the barber in a certain village shaves all and only those men of the village who do not shave themselves. Prove that the barber is a woman.
2. Prove that, in the set theory of Zermelo–Fraenkel, there does not exist a ‘universal’ set y , such that $\forall x (x \in y \Leftrightarrow x = x)$.
3. Prove that the impredicatively defined set α above is a real number according to Dedekind's definition as given in the text.

Intuitionistic Propositional Calculus

It seems that, over many centuries, no philosopher or mathematician ever seriously questioned Aristotle's *law of the excluded third*: for every proposition p , either p or not p , symbolically $p \vee \neg p$. In retrospect, it appears that Aristotle himself had some doubts about applying this law when talking about events in time, e.g., when p was the proposition: there will be a sea battle tomorrow. But in mathematics, which deals with unchanging entities, the law of the excluded third was accepted as gospel truth, as was the equivalent assertion: for every proposition p , $\neg\neg p \Rightarrow p$; two negations make an affirmation.

It was the topologist Luitzen Egbertus Jan Brouwer (1882–1966) who observed that all nonconstructive arguments in mathematics depend on Aristotle's law and he proposed that we simply drop it, together with all its consequences, at least when talking about infinite collections. Surprisingly, it turns out that, if one follows Brouwer's suggestion, one is still left with a rich logical system adequate for all constructive mathematics. We shall present the outlines of such a system here, starting with the propositional calculus.

We consider the logical symbols \top , \perp , \wedge , \vee and \Rightarrow , the first two counting as formulas, the last two being binary connectives between formulas, so that $A \wedge B$, $A \vee B$, and $A \Rightarrow B$ are formulas if A and B are. The usual reading of these formulas is as follows:

$$\begin{aligned}\top &\equiv \text{true,} \\ \perp &\equiv \text{false,} \\ A \wedge B &\equiv \text{(both) } A \text{ and } B,\end{aligned}$$

$$\begin{aligned} A \vee B &\equiv (\text{either}) A \text{ or } B, \\ A \Rightarrow B &\equiv \text{if } A \text{ then } B. \end{aligned}$$

However, *intuitionists*, as the followers of Brouwer are called, understand these words in a way subtly different from that of classical mathematicians, as we hope to make clear in the next chapter.

It is customary to make use of the *entailment* symbol \vdash , where

$$A_1, \dots, A_n \vdash B$$

means that the assumptions A_1, \dots, A_n entail the conclusion B , or that B may be deduced from A_1, \dots, A_n , n being any natural number. We often denote strings of formulas by capital Greek letters; thus $\Gamma \vdash B$ means that B may be inferred from the formulas in Γ . Note that Γ may be empty ($n = 0$), or consist of a single formula ($n = 1$).

We shall adopt the following *axioms*:

$$\vdash \top; \quad \perp \vdash A; \quad A \wedge B \vdash A; \quad A \wedge B \vdash B;$$

$$A, B \vdash A \wedge B; \quad A \vdash A \vee B; \quad B \vdash A \vee B; \quad A, A \Rightarrow B \vdash B.$$

The last of these has the Latin name ‘modus ponens’, which we shall abbreviate as MP. We also adopt two *rules of inference*:

$$\frac{\Gamma, A \vdash C, \quad \Gamma, B \vdash C}{\Gamma, A \vee B \vdash C}; \quad \frac{\Gamma, A \vdash B}{\Gamma \vdash A \Rightarrow B}.$$

These are called ‘argument by cases’, abbreviated AC, and ‘deduction rule’, abbreviated DR, respectively.

In classical logic, one would be quite happy to establish these axioms and rules of inference with the help of truth tables. In intuitionistic logic, we are not allowed to use truth tables, as they would also establish Aristotle’s $A \vee \neg A$ and $\neg \neg A \Rightarrow A$, which we have discarded, provided we define $\neg A \equiv A \Rightarrow \perp$, as we shall do from now on.

In addition to the above axioms and rules of inference, which describe the logical connectives, we also have the following, which describe the entailment symbol:

$A \vdash A$ (inferring a formula from itself);

$$\frac{\Lambda \vdash A \quad \Gamma, A \vdash B}{\Gamma, \Lambda \vdash B}$$
 (replacing an assumption by others which entail it);

$$\frac{\Gamma, A, B, \Delta \vdash C}{\Gamma, B, A, \Delta \vdash C}$$
 (interchanging two assumptions);

$$\frac{\Gamma \vdash B}{\Gamma, A \vdash B}$$
 (introducing a superfluous assumption);

$$\frac{\Gamma, A, A \vdash B}{\Gamma, A \vdash B}$$
 (contracting two identical assumptions into one).

These so-called *structural rules* were formally introduced by Gerhard Gentzen (1919–1945). Except for the axiom $A \vdash A$, they are often tacitly understood and not mentioned in actual arguments.

We will show how to establish deductions of the form $\Gamma \vdash B$ by looking at a few examples.

EXAMPLE 1. To prove that $A \vee B, \neg B \vdash A$.

Informally, we would argue as follows. We are given the assumptions $A \vee B$ and $B \Rightarrow \perp$. Suppose A . Then surely A , by the axiom $A \vdash A$. Suppose B . Then \perp by modus ponens from the second given assumption. Therefore A , by the axiom $\perp \vdash A$. Since A in either case, we invoke the argument by cases and infer that A holds in view of the given assumptions.

It is customary to rewrite such an argument more formally in a vertical fashion:

1	(1)	$A \vee B$	given
2	(2)	$B \Rightarrow \perp$	given
3	(3)	A	assumed
2,3	(4)	A	introducing a superfluous hypothesis
5	(5)	B	assumed
2,5	(6)	\perp	MP 2,5
2,5	(7)	A	by axiom $\perp \vdash A$
2,1	(8)	A	AC 4,7 replacing 3 and 5 by 1
1,2	(9)	A	interchanging two arguments

Note that the middle column contains the formulas given, assumed or inferred at different stages of the argument, numbered consecutively; the left column lists the numbers of all the hypotheses, given or assumed, upon which the formula in the middle column depends, and the right column indicates the justification for writing it down. The first two entries in the last line say precisely that $A \vee B, B \Rightarrow \perp \vdash A$, as was to be proved.

EXAMPLE 2. To prove that $A \vdash \neg\neg A$.

Here is the informal argument: we are given A . Suppose $A \Rightarrow \perp$. Then \perp by modus ponens. Therefore, $(A \Rightarrow \perp) \Rightarrow \perp$ by the deduction rule.

Formally:

1	(1)	A	given
2	(2)	$A \Rightarrow \perp$	assumed
1,2	(3)	\perp	MP 1,2
1	(4)	$(A \Rightarrow \perp) \Rightarrow \perp$	DR 2,3

EXAMPLE 3. To prove that $\vdash A \Rightarrow (B \Rightarrow A)$.

We shall give the formal argument only.

1	(1)	A	assumed
2	(2)	B	assumed
1,2	(3)	A	introducing a superfluous hypothesis
1	(4)	$B \Rightarrow A$	DR 2,3
	(5)	$A \Rightarrow (B \Rightarrow A)$	DR 1,4

After some practice, the student may stop mentioning the structural rules, such as line (3) above, or lines (4) and (9) in Example 1.

EXAMPLE 4. To prove that $A \Rightarrow B \vdash \neg B \Rightarrow \neg A$.

1	(1)	$A \Rightarrow B$	given
2	(2)	$B \Rightarrow \perp$	assumed
3	(3)	A	assumed
1,3	(4)	B	MP 1,3
1,2,3	(5)	\perp	MP 2,4
1,2	(6)	$A \Rightarrow \perp$	DR 3,5
1	(7)	$\neg B \Rightarrow \neg A$	DR 2,6

EXAMPLE 5. To prove that $\vdash \neg\neg(A \vee A)$.

1	(1)	$(A \vee \neg A) \Rightarrow \perp$	assumed
2	(2)	A	assumed
2	(3)	$A \vee \neg A$	axiom $A \vdash A \vee B$
1,2	(4)	\perp	MP 1,3
1	(5)	$A \Rightarrow \perp$	DR 2,4
1	(6)	$A \vee \neg A$	axiom $B \vdash A \vee B$
1	(7)	\perp	MP 1,6
	(8)	$((A \vee \neg A) \Rightarrow \perp) \Rightarrow \perp$	DR 1,7

Exercises

Prove the following.

1. $A \vee \neg A \vdash \neg\neg A \Rightarrow A$.
2. $((A \Rightarrow B) \Rightarrow A) \vdash \neg\neg A$.
3. $C \Rightarrow (A \wedge B) \vdash (C \Rightarrow A) \wedge (C \Rightarrow B)$.
4. $(C \Rightarrow A) \wedge (C \Rightarrow B) \vdash C \Rightarrow (A \wedge B)$.
5. $A \Rightarrow (B \Rightarrow C) \vdash (A \wedge B) \Rightarrow C$.
6. $(A \wedge B) \Rightarrow C \vdash A \Rightarrow (B \Rightarrow C)$.

Note that the classical result $((A \Rightarrow B) \Rightarrow A) \vdash A$ follows from (2) if $\neg\neg A \vdash A$, but it does not hold intuitionistically.

How to Interpret Intuitionistic Logic

To explain the subtle difference between intuitionistic and classical logic, we shall present an intuitionistic interpretation of the logical connectives that goes back to Brouwer, Heyting and Kolmogorov. It involves talking about *reasons* for a formula. A *reason* for A may be thought of as a proof of A from some suitable assumption. We would like to say

- there is exactly one reason for \top (namely, quoting the axiom);
- there is no reason for \perp ;
- a reason for $A \wedge B$ consists of a reason for A and a reason for B ;
- a reason for $A \vee B$ is a reason for A or a reason for B ;
- a reason for $B \Rightarrow C$ is a rule for converting a reason for B into a reason for C .

Now compare these statements with the following statements about sets (see Chapter 13, where 0 was defined as the empty set and 1 as $\{0\}$):

- there is exactly one element of 1 ;
- there is no element of 0 ;
- an element of $A \times B$ is a pair of elements of A and B , respectively;
- an element of $A + B$ is an element of A or an element of B ;

- an element of C^B is a function that converts an element of B into an element of C .

Comparing intuitionistic logic with the arithmetic of sets, we are led to the following analogies:

$$\begin{array}{ll}
 \top & 1 \\
 \perp & 0 \\
 A \wedge B & A \times B \\
 A \vee B & A + B \\
 B \Rightarrow C & C^B
 \end{array}$$

Moreover, a deduction from A to B , namely, an argument showing that $A \vdash B$, corresponds to a mapping $A \rightarrow B$. If there is a deduction from A to B and a deduction from B to A , we shall write $A \vdash\vdash B$. This corresponds to mappings $A \rightarrow B$ and $B \rightarrow A$ and we may write $A \leftrightarrow B$. Frequently these two mappings are inverse to one another, so we have a one-to-one correspondence between A and B , that is, $A \cong B$.

For example, we can prove intuitionistically that

$$\begin{aligned}
 C \Rightarrow (A \wedge B) &\vdash\vdash (C \Rightarrow A) \wedge (C \Rightarrow B), \\
 A \Rightarrow (B \Rightarrow C) &\vdash\vdash (A \wedge B) \Rightarrow C, \\
 (A \vee B) \Rightarrow C &\vdash\vdash (A \Rightarrow C) \wedge (B \Rightarrow C).
 \end{aligned}$$

These equivalences correspond to the following one-to-one mappings between sets:

$$\begin{aligned}
 (A \times B)^C &\cong A^C \times B^C, \\
 (C^B)^A &\cong C^{A \times B}, \\
 C^{A+B} &\cong C^A \times C^B.
 \end{aligned}$$

As we saw in Chapter 13, these are just generalizations of the familiar laws of arithmetic:

$$\begin{aligned}
 (a \times b)^c &= a^c \times b^c, \\
 (c^b)^a &= c^{a \times b}, \\
 c^{a+b} &= c^a \times c^b.
 \end{aligned}$$

One cannot but be impressed by the remarkable unity pervading logic, set theory and arithmetic.

A word of warning: $A \leftrightarrow B$ does not always mean $A \cong B$. For example, the intuitionistic equivalence

$$A \Rightarrow (B \Rightarrow A) \vdash\vdash \top$$

translates into

$$(A^B)^A \leftrightarrow 1,$$

with mappings in both directions; but the 1-to-1 correspondence $(A^B)^A \cong 1$ holds only if $A \cong 1$ or $A \cong 0$ or $B \cong 0$.

We should also point out that the interpretation of $A \vee B$ advocated here does not work for classical logic. As we see in Chapter 28, there is a formula G in number theory for which one can prove neither G nor $\neg G$. According to the intended interpretation of $A \vee B$, $G \vee \neg G$ has no proof. Yet, classically, $G \vee \neg G$ is just a case of Aristotle's axiom of the excluded third and thus quoting this axiom would constitute a proof.

Exercises

1. Take any of the one-to-one correspondences of Chapter 13, translate it into a statement of intuitionistic logic and prove the latter.
2. Take any intuitionistic theorem of the form $A \vdash B$ and find the corresponding mapping $A \rightarrow B$ between sets.

Intuitionistic Predicate Calculus

We shall be dealing with formulas A which contain so-called *free variables* such as x in $x^2 + 2 = 0$ or x and y in $xy + x = y$. To indicate which variables may be present we often write $A(x)$ or $A(x, y)$ instead of just A .

From $A(x)$ one may obtain the formulas $\forall_x A(x)$ and $\exists_x A(x)$ in which x is no longer free; it is *bound* to the *universal quantifier* \forall_x , meaning *for all* x , or to the *existential quantifier* \exists_x , meaning *for some* x . Similarly, one may form $\forall_x \forall_y A(x, y)$, $\forall_x \exists_y A(x, y)$, etc.

From $A(x)$ one may also obtain the formula $A(t)$, the result of substituting the *term* t for x . Here t may be 5 or y or $x + 3y$ or whatever — it may even be x . However, if $A(x)$ is $\exists_y B(x, y)$, we are not supposed to substitute y or $x + 3y$ for x , because y is no longer free in $\exists_y B(y, y)$ or $\exists_y B(3x + y, y)$. This prohibition is spelled out in the following definition: t is *substitutable* for x in $A(x)$ provided any free variable in t (maybe t itself) remains free in $A(t)$. This definition is needed for stating the following two axioms:

$$\forall_x A(x) \vdash A(t) \quad (\text{universal specification}),$$

$$A(t) \vdash \exists_x A(x) \quad (\text{existential generalization}),$$

— subject to the restriction that t is substitutable for x in $A(x)$.

Were it not for this restriction, we would have as a special case of the first axiom

$$\forall_x \exists_y B(x, y) \vdash \exists_y B(y, y)$$

and from ‘everybody blames somebody’ we could infer that ‘somebody blames himself’.

In addition to the above two axioms, we shall also adopt the following two rules of inference:

$$\frac{\Gamma \vdash A(x)}{\Gamma \vdash \forall_x A(x)} \quad (\text{universal generalization}),$$

provided x is not free in Γ ;

$$\frac{\Gamma, A(x) \vdash B}{\Gamma, \exists_x A(x) \vdash B} \quad (\text{existential specification}),$$

provided x is not free in Γ or B .

We abbreviate the names of these axioms and rules of inference by US, EG, UG, and ES, respectively. The reason for the restriction on UG, for example, is to avoid inferring from ‘ x is afraid’ that ‘everyone is afraid’.

We shall present some examples to illustrate arguments involving quantifiers.

EXAMPLE 1. To prove $\forall_x (F(x) \wedge G(x)) \vdash \forall_x F(x) \wedge \forall_x G(x)$.

1		(1)	$\forall_x (F(x) \wedge G(x))$	given
1		(2)	$F(x) \wedge G(x)$	US 1
1		(3)	$F(x)$	axiom for \wedge , 2
1		(4)	$\forall_x F(x)$	UG 3 (x not free in 1)
1		(5)	$G(x)$	axiom for \wedge , 2
1		(6)	$\forall_x G(x)$	UG 5 (x not free in 1)
1		(7)	$\forall_x F(x) \wedge \forall_x G(x)$	axiom for \wedge , 4,6

EXAMPLE 2. To prove $\exists_y \forall_x F(x, y) \vdash \forall_x \exists_y F(x, y)$.

1		(1)	$\exists_y \forall_x F(x, y)$	given
2		(2)	$\forall_x F(x, y)$	assumed
2		(3)	$F(x, y)$	US 2
2		(4)	$\exists_y F(x, y)$	EG 3
1		(5)	$\exists_y \forall_x F(x, y)$	ES 4 (y not free in 4)
1		(6)	$\forall_x \exists_y F(x, y)$	UG 5 (x not free in 1)

EXAMPLE 3. To prove $\neg\exists_x F(x) \vdash \forall_x \neg F(x)$.

1	(1)	$\exists_x F(x) \Rightarrow \perp$	given
2	(2)	$F(x)$	assumed
2	(3)	$\exists_x F(x)$	EG 2
1,2	(4)	\perp	MP 1,2
1	(5)	$F(x) \Rightarrow \perp$	DR 4,2
1	(6)	$\forall_x (F(x) \Rightarrow \perp)$	UG 5 (x not free in 1)

Exercises

1. Prove the converse of Example 1.
2. What goes wrong if you try to prove the converse of Example 2?
3. Can you prove the converse of Example 3?
4. Prove that $\neg\forall_x \neg F(x) \vdash \neg\neg\exists_x F(x)$, using Example 3 above and Example 4 of Chapter 24. Classically, but not intuitionistically, one can infer from this that $\neg\forall_x \neg F(x) \vdash \exists_x F(x)$. If we could prove $\neg\forall_x \neg F(x)$, we would have a nonconstructive proof of $\exists_x F(x)$.
5. Prove that $\exists_x \forall_y F(x, y) \vdash \exists_x F(x, x)$.

Intuitionistic Type Theory

This chapter is an adaptation of the appendix in Couture and Lambek [1991], giving a brief overview of a recent formulation of type theory in Lambek and Scott [1986], which is adequate for elementary mathematics, including arithmetic and analysis, when treated constructively. As far as we know, the only proofs in these disciplines which are essentially nonconstructive depend on the *axiom of choice*. One formulation of this axiom asserts that, for any nonempty collection of nonempty sets, there exists a set containing exactly one element from each of the given sets.

From basic types 1 , Ω and N one builds others by two processes: if A is a type so is PA ; if A and B are types, so is $A \times B$. The intended meaning of these types is as follows:

- 1 is the type of a specified single entity (introduced for convenience);
- Ω is the type of truth values or propositions (here there are more than two truth values);
- N is the type of natural numbers;
- PA is the type of sets of entities of type A ;
- $A \times B$ is the type of pairs of entities of types A and B , respectively.

We allow arbitrarily many variables of each type and write $x \in A$ to mean that x is a variable of type A . In addition, we construct terms of different types inductively as follows:

$$\frac{1 \quad \Omega \quad N \quad PA \quad A \times B}{* \quad a = a' \quad 0 \quad \{x \in A \mid \phi(x)\} \quad < a, b >} \\ a \in \alpha \quad Sn$$

— it being assumed that a and a' are terms of type A already constructed, α of type PA , n of type N , $\phi(x)$ of type Ω and b of type B .

Logical symbols may be defined as follows:

$$\begin{aligned} \top &\equiv * = *, \\ p \vee q &\equiv < p, q > = < \top, \top >, \\ p \Rightarrow q &\equiv p \wedge q = p, \\ \forall_{x \in A} \phi(x) &\equiv \{x \in A \mid \phi(x)\} = \{x \in A \mid \top\}, \end{aligned}$$

where it is understood that p, q and $\phi(x)$ are terms of type Ω . From these symbols one may define others, taking care not to make implicit use of De Morgan's rules (Prawitz [1965]):

$$\begin{aligned} \perp &\equiv \forall_{t \in \Omega} t, \\ \neg p &\equiv \forall_{t \in \Omega} (p \Rightarrow t), \\ p \wedge q &\equiv \forall_{t \in \Omega} ((p \Rightarrow t) \wedge (q \Rightarrow t)) \Rightarrow t, \\ \exists_{x \in A} \phi(x) &\equiv \forall_{t \in \Omega} ((\forall_{x \in A} (\phi(x) \Rightarrow t)) \Rightarrow t). \end{aligned}$$

Other symbols, such as appear in $\exists_{x \in A} \phi(x)$, $\{a\}$, $\alpha \subseteq \beta$, $\alpha \times \beta$, etc., are defined in the usual fashion.

Axioms and rules of inference are stated in terms of a deduction symbol \vdash_X , where X is a finite set of variables. The permissible deductions take the form

$$p_1, \dots, p_n \vdash_X p_{n+1},$$

where the p_i are terms of type Ω and X contains all the variables which occur freely in the p_i . (When X is empty, the subscript may be omitted.) The axioms and rules of inference hold no surprises. For the purpose of illustration, here are a few special cases:

- $p \vdash p$;
- $\frac{\phi(x) \vdash_{\{x\}} \psi(x)}{\phi(\alpha) \vdash \psi(\alpha)}$;
- $< a, b > = < c, d > \vdash a = c$;
- $\frac{\phi(x) \vdash_{\{x\}} \phi(Sx)}{\phi(0) \vdash_{\{x\}} \phi(x)}$.

The reader may recognize the last mentioned rule of inference as the principle of mathematical induction. Although we have not stated all the axioms

and rules of inference here, they will of course imply the usual axioms and rules for the intuitionistic propositional and predicate calculi.

For a less formal treatment of constructive arithmetic and analysis, the reader may consult Goodstein [1970] and Bishop [1967], respectively.

Gödel's Theorems

A formal language is (among other things) a system for dealing with strings of symbols. An interpretation of these symbols is called a *model*. For example, the Lambda Calculus tells us how to arrange marks such as ‘ and λ_x . The interpretation of these marks in terms of functional application and functional abstraction is a model.

The formal languages that interest us here contain the notion of a proof. A *proof* is a finite sequence of formulas, each of which is either an axiom or follows from some previous members of the sequence by a rule of inference. We call a formal language *consistent* if there is no proof in the language whose last line is \perp , that is, it does not contain the proof of a contradiction.

In this chapter we consider a formal language L which is adequate for *arithmetic*. We assume that L includes the intuitionistic predicate calculus together with some axioms for arithmetic, and we assume that all the basic laws of arithmetic are provable in L . For example, L might be the type theory considered in the previous chapter.

In 1930, Kurt Gödel (1906–1978) proved a completeness result for the classical predicate calculus. In 1950, this result was extended by Leon Henkin to classical type theory. It was later extended to intuitionistic type theory (Lambek and Scott [1986]). This completeness result may be expressed as follows:

Theorem 28.1. (Completeness)

A formula is provable in L if it is true under all possible interpretations of the nonlogical symbols in L , i.e., in all models of L .

We shall not give a proof here, but we remark that the proof of this theorem depends on the axiom of choice and is thus unacceptable to intuitionists. The converse of this theorem is called the soundness theorem; its proof is straightforward and is acceptable to intuitionists. As a corollary of the completeness-soundness result, L is consistent if and only if it has a model.

If n is a natural number, let $S^n 0$ be the expression in L formed by placing the letter S n times before the symbol 0 . This expression will normally be interpreted as the natural number n . We call a model ω -complete if, for any formula $A(x)$ of L , x being of type N , whenever $A(S^n 0)$ is true in that model for each natural number n , then $\forall_{x \in N} A(x)$ is also true in that model.

In 1931, Gödel proved his *incompleteness theorem* for arithmetic, which may be expressed in our terminology as follows:

Theorem 28.2. (Incompleteness)

There is a formula in L which is true in any ω -complete model, but not provable in L , assuming L to be consistent.

Combining this with the completeness theorem, we may conclude that some models are not ω -complete.

To an intuitionist, the notion of truth is equivalent to that of knowability, which we shall here interpret as provability. Thus we may conclude that the world in which we live is ω -complete provided, whenever $A(S^n 0)$ has a proof for each n , then $\forall_{x \in N} A(x)$ does also. However, there is no particular reason to believe that this is the case. Even if, for each n , the formula $A(S^n 0)$ showed up as the last line of a proof, it would not guarantee that $\forall_{x \in N} A(x)$ showed up as the last line of a proof. Hence the intuitionist has no particular reason to think of the world we live in as ω -complete.

Platonists, of whom Gödel was one, see truth as the property of an eternal and immutable reality which is independent of finite human minds. A classical Platonist believes that the real world contains an infinite collection of natural numbers. Now if it is true of each of these numbers that it has a property A , then it is true that they all have property A , whence the real world is an ω -complete model of L . Since Gödel was a classical Platonist, he concluded from his incompleteness theorem that there is a formula of arithmetic which is eternally true, but which is not the last line of a proof in L . (Note that for Platonists, the language of arithmetic is about the real world, hence it is consistent.)

Proof of Gödel's Incompleteness Theorem

We shall now sketch a proof of the incompleteness result in a manner acceptable to intuitionists. We begin with a lemma.

Lemma 29.1. (Gödel's Lemma)

Suppose $R(m, n)$ is a recursive relation between the natural numbers m and n . That is, assume that, for any two numbers m and n , there is a finite effective procedure for deciding whether they are in the relation R . Then there is a formula $F(x, y)$ in L , with x and y of type N , such that

- *if $R(m, n)$ then $\vdash F(S^m 0, S^n 0)$,*
- *if not $R(m, n)$ then $\vdash \neg F(S^m 0, S^n 0)$.*

(Here \vdash means 'there is a proof in L that'.)

As an example of this lemma, let R be the relation 'is 1 greater than'. Let $F(x, y)$ be the formula $x = Sy$. Then, if m is 1 greater than n , it is provable in L that $S^m 0 = SS^n 0$, and if m is not greater than n , it is provable in L that $S^m 0 \neq SS^n 0$.

We shall not prove this lemma here, but it should seem reasonable, inasmuch as L is meant to capture ordinary number theory.

Theorem 29.2.

(Gödel's Incompleteness Theorem (Semantic Version))

If L is consistent, there is a formula in L which is true in any ω -complete model but not provable in L .

Proof: We begin by enumerating all expressions in L of type PN , that is, of the form $\{x \in N | A(x)\}$, where $A(x)$ is a formula in L , with x of type

N . Since there are only countably many finite strings of symbols, there are only countably many expressions of this form. Call them $\alpha_1, \alpha_2, \dots$

We next enumerate all the proofs of L . There are only countably many — assuming that proofs are finite — since there are only countably many finite sequences of formulas. Call the proofs P_1, P_2, \dots

Let $R(m, n)$ be the relation which is satisfied just in case P_n is a proof of the formula $S^m 0 \in \alpha_m$. $R(m, n)$ is recursive, since we can decide whether it is true for a given m and n by looking to see whether $S^m 0 \in \alpha_m$ is the last line of the proof P_n .

By Gödel's Lemma, there is a formula $F(x, y)$ such that, if $R(m, n)$, then $\vdash F(S^m 0, S^n 0)$ and, if not $R(m, n)$, then $\vdash \neg F(S^m 0, S^n 0)$. Consider $\{x \in N \mid \neg \exists y \in N F(x, y)\}$. This is one of the $\alpha_1, \alpha_2, \dots$, say α_g . Then, from the definition of α_g , it is provable in L that

$$S^g 0 \in \alpha_g \iff \neg \exists y \in N F(S^g 0, y). \quad (\clubsuit)$$

In both intuitionistic and ordinary logic, it follows that

$$\vdash \neg S^g 0 \in \alpha_g \iff \vdash \exists y \in N F(S^g 0, y).$$

Consider the formula $S^g 0 \in \alpha_g$. If it is provable in L , then there is some natural number n such that P_n is a proof of $S^g 0 \in \alpha_g$. In that case, $R(g, n)$. But then $\vdash F(S^g 0, S^n 0)$, and hence $\vdash \exists y \in N F(S^g 0, y)$. We thus have $\vdash \neg S^g 0 \in \alpha_g$.

Hence it is not the case that $S^g 0 \in \alpha_g$ (assuming L is consistent). That is, for all n , not $R(g, n)$, and hence, for all n , $\vdash \neg F(S^g 0, S^n 0)$ (by Gödel's Lemma). Call this result $(*)$, for future reference.

From the soundness result, we may conclude that, for all n , $\neg F(S^g 0, S^n 0)$ is true in all models of L . Hence in any ω -complete model of L , the proposition $\forall y \in N \neg F(S^g 0, y)$ is true.

Even intuitionistic logic allows us to infer $\neg \exists y G(y)$ from $\forall y \neg G(y)$. Hence we may conclude that $\neg \exists y \in N F(S^g 0, y)$ is true in any ω -complete model, and thus $S^g 0 \in \alpha_g$ is true in any such model by (\clubsuit) . Thus, although $S^g 0 \in \alpha_g$ is not provable in L (since L is consistent), $S^g 0 \in \alpha_g$ is nonetheless true in all ω -complete models.

Exercises

1. Prove the two claims made about inferences in intuitionistic logic.
2. If $R(m, n)$ means ' m and n both equal 3', find the corresponding formula $F(x, y)$ whose existence is asserted by Gödel's Lemma.

More about Gödel's Theorems

We say that a language L is ω -consistent provided, for any formula $A(x)$ of L , if $A(S^n 0)$ is provable for each natural number n , then it is not the case that $\vdash \neg \forall_{y \in N} A(y)$. It is not hard to show that ω -consistency implies consistency.

If L is ω -consistent, then it follows from $(*)$ in the previous chapter that it is not the case that

$$\vdash \neg \forall_{y \in N} \neg F(S^g 0, y).$$

From (\clubsuit) in the previous chapter, it now follows that it is not the case that $\vdash \neg S^g 0 \in a_g$. This gives us

Theorem 30.1.

(Gödel's Incompleteness Theorem (Syntactic Version))

If L is ω -consistent, there is a formula G such that neither G nor $\neg G$ is provable in L .

Rosser showed that ω -consistency here can be replaced by plain consistency. His proof is short, but tricky, and we shall skip it.

If neither a statement nor its negation is provable in a language, it is called *undecidable* relative to that language. There is no way to get rid of all undecidable statements in a language by adjoining a finite number of new axioms. If we added $S^g 0 \in a_g$ to the axioms of L , then there would still be a g' such that $S^{g'} 0 \in a_{g'}$ is an undecidable statement relative to this new language.

Hilbert's second problem was to prove the consistency of arithmetic using

only formal arithmetic to do so. Gödel's Incompleteness Theorem implies that this is impossible. For let *Cons* be a statement in L which expresses the idea that there is no n such that P_n is a proof of \perp . Then, if L is consistent, and *if there is a proof in L of this*, we shall have $\vdash \text{Cons}$. If L is not consistent, we can prove anything, so we can have $\vdash \neg \text{Cons}$ as well as $\vdash \text{Cons}$.

One can formalize the proof of Gödel's Incompleteness Theorem to show that $\vdash \text{Cons} \Rightarrow S^g 0 \in a_g$. Hence, if there were a proof in L that L was consistent, namely, if $\vdash \text{Cons}$, then we would have $\vdash S^g 0 \in a_g$. But this we have seen is not the case.

Gödel's Incompleteness Theorem, being a metamathematical result, has different implications depending on one's conception of mathematics. A classical *formalist* views arithmetic as nothing more than strings of symbols, manipulated according to certain rules. He does not want to rely on any interpretation of these symbols in order to ensure the consistency of the rules. He would hope that consistency could be established within the system itself. Gödel's Theorem shows, against the hopes of the classical formalist, that the consistency of arithmetic cannot be demonstrated within arithmetic. Moreover, it seems that the idea of truth cannot be captured by the notion of provability.

From an *intuitionist* point of view, Gödel's Incompleteness Theorem is not unwelcome. For intuitionists, 'true' should mean 'provable' (with a finite proof). They maintain that there are some statements which are neither true nor false; Gödel's result merely confirms this belief by showing that there are statements which are neither provable nor disprovable.

For the classical *Platonist*, the Incompleteness Theorem shows that there are statements true in the real world which are not provable. In other words, the realities of mathematics are too profound to be captured by any finite axiom system. There will always be truths in mathematics which cannot be cranked out by a computer, but which must await new philosophical insights for their discovery.

Concrete Categories

In the 20th century, we find a great deal of concrete, practical mathematics. Statistics is flourishing, the computer has proved the Four Colour Theorem, and numbers with upwards of two hundred digits can be factored.

Another trend in the 20th century is a degree of abstraction never seen before in mathematics. For example, the study of the Euclidean plane has been replaced by the study of vector spaces and topological spaces that abstract some of its properties. A prominent and influential proponent of this trend in algebra was Emmy Noether (1882–1935). The supreme abstraction is the notion of a *category*, to which we shall turn our attention in this chapter.

Nowadays, when studying vector spaces, we are forced to look also at linear transformations. Similarly, when studying topological spaces, we are led to continuous mappings. When studying groups (themselves an abstraction of permutation groups), we have to look at homomorphisms. To abstract the properties which these examples have in common, we introduce the notion of a ‘concrete category’.

A *concrete category* is a class of sets, each endowed with a certain structure, together with the class of all functions which map one set to another while preserving this structure.

EXAMPLE 31.1

The class of sets together with the class of all functions between them is a concrete category. Here there is no structure to preserve, so the condition on the functions is trivially satisfied.

EXAMPLE 31.2

A *monoid* is a set containing a special *identity* element 1 together with a binary operation \cdot between members of that set, such that $(a \cdot b) \cdot c = a \cdot (b \cdot c)$, and $1 \cdot a = a \cdot 1 = a$. A *monoid homomorphism* is a function f from a monoid A to another monoid A' which preserves structure: $f(a \cdot b) = f(a) \cdot f(b)$ and $f(1) = 1$. (When we write ' $f(1) = 1$ ' it is understood that the first 1 is the special element of A and the second 1 is the special element of A' .) For example, taking 0 as the special element and $+$ as the binary operation, the natural numbers form a monoid. As another example, the singleton set $\{1\}$ is a monoid with $1 \cdot 1 = 1$. The mapping from the natural numbers to $\{1\}$ is a monoid homomorphism. Note that group homomorphisms are a special case of monoid homomorphisms. The class of monoids, together with the monoid homomorphisms forms a concrete category.

EXAMPLE 31.3

A *pre-ordered set* is a set, together with a binary relation \leq on that set which is reflexive and transitive: $a \leq a$, and if $a \leq b$ and $b \leq c$ then $a \leq c$. A *monotone mapping* is a function f from a pre-ordered set A to a pre-ordered set A' that preserves the order: if $a \leq b$ then $f(a) \leq f(b)$.

For example, both the natural numbers and the even numbers together with the relation \leq are pre-ordered sets. The doubling function maps the natural numbers into the set of even numbers in an order preserving way, and is thus a monotone mapping. The collection of all pre-ordered sets together with all monotone mappings between them forms a concrete category.

Note that in all cases the identity function on a set will preserve its structure. Also, the composition of two structure preserving functions will preserve structure.

More abstract than a concrete category is a *category*, which we shall define in Chapter 32.

Exercise

1. Verify that groups and group homomorphisms form a concrete category.

Graphs and Categories

A *graph* (more precisely: an *oriented multigraph*) consists of a class of *arrows* (or *directed edges*) together with a class of *objects* (or *nodes*), and also two mappings from the class of arrows to the class of objects. The mappings are called S (*source* or *domain*) and T (*target* or *codomain*).

$$\begin{array}{ccc} \{\text{arrows}\} & \begin{array}{c} \xrightarrow{\text{source}} \\ \xrightarrow{\hspace{1cm}} \\ \xrightarrow{\text{target}} \end{array} & \{\text{objects}\} \end{array}$$

If f is an arrow, $S(f) = A$ and $T(f) = B$, we write

$$f : A \rightarrow B \text{ or } A \xrightarrow{f} B.$$

A *category* is a graph subject to the following conditions:

1. associated with any two arrows $f : A \rightarrow B$ and $g : B \rightarrow C$ (so that $T(f) = S(g)$) is an arrow $g \circ f : A \rightarrow C$;
2. if $f : A \rightarrow B$, $g : B \rightarrow C$ and $h : C \rightarrow D$, then $(h \circ g) \circ f = h \circ (g \circ f)$;
3. associated with each object A , there is an identity arrow 1_A , whose source and target are A ;
4. if $f : A \rightarrow B$ then $f \circ 1_A = f$; if $g : B \rightarrow A$ then $1_A \circ g = g$.

A concrete category is a category whose objects are sets with structure and whose arrows are the structure preserving functions between them. In concrete categories, 1_A is the identity map on A , and \circ is function composition.

EXAMPLE 32.1

Let A be any set. Let the class of objects be A . Let the class of arrows be A . Let $S(a) = T(a) = a$ for all $a \in A$. Let 1_a be a and let $a \circ a = a$. Then the conditions for a category are satisfied and we have what is called the *discrete* category corresponding to A . So a set may be viewed as a category.

EXAMPLE 32.2

Let $(A, 1, \cdot)$ be a monoid (with special object 1 and binary operation \cdot). For the class of objects take the singleton set $\{*\}$. For the class of arrows, take A . Let $S(a) = T(a) = *$ for all $a \in A$. Let 1_* be 1 (the monoid identity) and let $a \circ b$ be $a \cdot b$. The graph we have just constructed is a category, thanks to the structure of the monoid. In this way a monoid may be viewed as a category. If we consider the monoid of the natural numbers with addition, in this way, as a category, then what is the number 2? It is the unique arrow in the natural number monoid, viewed as a category, which can be written as a composition of nonidentity arrows in exactly one way.

EXAMPLE 32.3

Let (A, \leq) be a pre-ordered set. Let the class of objects be A , and let the class of arrows be $\{(a, b) | a \leq b\}$. (Here a and b are assumed to be elements of A .) Let $S((a, b)) = a$ and $T((a, b)) = b$. By the reflexivity of \leq , (a, a) is an arrow for all $a \in A$. Let this be the identity arrow associated with a , so that $1_a = (a, a)$. Define the composition of arrows thus: $(b, c) \circ (a, b) = (a, c)$. By transitivity of \leq , we know that (a, c) is an arrow if (a, b) and (b, c) are arrows. Again, we have a category.

The examples from Chapter 31 show that *many interesting objects in mathematics congregate in categories*. The examples in this chapter illustrate that *interesting mathematical entities may often be viewed as categories*.

Exercise

1. Show that there is a category with exactly two objects and exactly one nonidentity arrow. (This category is sometimes said to be the number 2.)

33

Functors

If A and B are two categories, a *functor* F from A to B is a mapping sending objects of A to objects of B and, at the same time, a mapping sending arrows of A to arrows of B , so that

1. if g is any arrow in A with source a and target a' , then $F(g)$ is an arrow in B with source $F(a)$ and target $F(a')$;
2. $F(1_a) = 1_{F(a)}$;
3. $F(g \circ h) = F(g) \circ F(h)$.

We saw in Chapter 32 that sets, monoids and pre-ordered sets may all be viewed as categories. How do the structure preserving mappings between such entities compare with the functors between them when they are viewed as categories? We answer the question as follows.

EXAMPLE 33.1

Suppose A and B are sets, each viewed as a category (Example 32.1). Let F be any map from A to B . If $a \in A$, then F maps a to $F(a)$. But a is just 1_a and $F(a)$ is just $1_{F(a)}$. Hence $F(1_a) = 1_{F(a)}$. If a and b are any elements of A , then they are also arrows of A . If $a \neq b$ then $S(a) = T(a) = a \neq b = S(b) = T(b)$, so they cannot be composed. If $a = b$ then $F(a \circ a) = F(a) = F(a) \circ F(a)$. Hence F is a functor.

EXAMPLE 33.2

Suppose A and B are monoids, each viewed as a category. Let F be a monoid homomorphism from A to B . Without loss of generality, we may suppose that $\{*\}$ is the class of objects for both categories A and B (see Example 32.2). Suppose $F(*) = *$. Since F is a homomorphism, $F(1_*) = 1_{F(*)}$. Moreover, $F(a \circ a') = F(a \cdot a') = F(a) \cdot F(a') = F(a) \circ F(a')$. Hence F is a functor.

EXAMPLE 33.3

A monotone mapping between pre-ordered sets may be viewed as a functor. The details are left as an exercise.

The next three examples illustrate the observation that *many entities of interest in mathematics may be viewed as functors*.

EXAMPLE 33.4

A set may be viewed as a category, as we saw in Example 32.1. It can also be viewed as a functor. Let A be the discrete one-element category and B the category of sets. For any set S there is a unique functor from A to B such that $F(1_A)$ is the identity function on S . This functor can be viewed as the set S .

EXAMPLE 33.5

Let A be a category with two objects, \mathbf{a} and \mathbf{o} , and with four arrows: $1_{\mathbf{a}}$, $1_{\mathbf{o}}$, $s : \mathbf{a} \rightarrow \mathbf{o}$, and $t : \mathbf{a} \rightarrow \mathbf{o}$. This category may be pictured thus:

$$\begin{array}{ccc} \mathbf{a} & \xrightarrow{\quad} & \mathbf{o} \\ & \xrightarrow{\quad} & \end{array}$$

Suppose F maps \mathbf{a} to a set X and \mathbf{o} to a set Y . Suppose $F(s)$ and $F(t)$ are functions with domain X and codomain Y . Then F is a functor from A to the category of sets. We can think of this functor as a graph with class of objects Y , class of arrows X , source mapping $F(s)$ and target mapping $F(t)$.

EXAMPLE 33.6

Let M be a monoid and X a set. Suppose $m : M \times X \rightarrow X$ is a function such that, for elements a and a' of M and b of X , $m(a \cdot a', b) = m(a, m(a', b))$. Suppose also that $m(1, b) = b$. Then (M, X, m) is an M -set. For example, M might be the monoid of positive integers with multiplication, and X might be the set of segments constructible in Euclidean geometry, and $m(a, CD)$ might be the function mapping a segment CD to

a segment CE (with E on CD produced) which is a times as long as CD . One often writes $m(a, b)$ as ab .

How can we make an M -set into a functor? Let F map the category M to the category of sets so that $F(*)$ is X and, for any $a \in M$, $F(a) : X \rightarrow X$ is the function which sends $b \in X$ onto $F(a)(b) = ab$. Thus $F(1)$ is the identity function on X and $F(a \cdot a') = F(a) \circ F(a')$. For, if $b \in X$ then $F(a \cdot a')(b) = (aa')b = a(a'b) = F(a)(F(a')(b))$.

Call a category *small* if its class of objects and its class of arrows are both *sets*. (The distinction between 'set' and 'class' is made clear in the set theory of Gödel and Bernays.) *Cat* is the category whose objects are small categories and whose arrows are the functors from one small category to another. The fact that the small categories themselves form a category again illustrates the slogan that 'interesting objects congregate in categories'.

Exercises

1. Complete the details of Example 33.3 on pre-ordered sets.
2. Show in detail that *Cat* is indeed a category.
3. A graph is *small* if its class of objects and its class of arrows are both sets. Show that there is a category *Grph* whose objects are small graphs, and whose arrows are like functors, except that they need not satisfy the equations (2) and (3) in the definition of a functor. (*Grph* will be discussed in some detail in the following section.)

Natural Transformations

Category theory began in 1945 with Eilenberg and Mac Lane's article 'General Theory of Natural Equivalences'. In this chapter we will investigate the notion of a 'natural equivalence'.

Let A and B be categories and let F and G be functors from A to B . A *natural transformation* t from F to G is a mapping that assigns to every object a of A an arrow $t(a)$ in B from $F(a)$ to $G(a)$, such that, for any arrow f in A from a to b , $G(f) \circ t(a) = t(b) \circ F(f)$. This can be pictured as follows:

$$\begin{array}{ccc}
 & t(a) & \\
 F(a) & \longrightarrow & G(a) \\
 F(f) \downarrow & & \downarrow G(f) \\
 & t(b) & \\
 F(b) & \longrightarrow & G(b)
 \end{array}$$

Note that a and b are objects in category A , whereas all of the objects and arrows in the picture are in B .

EXAMPLE 34.1

We saw in Example 33.4 that a set S can be viewed as a functor F from the discrete one-object category to the category of sets, in such a way that F takes its object to S and its arrow to the identity function on S . Let S and S' be two sets, and F and F' the corresponding functors. Given a function f from S to S' , let t be a mapping that assigns to the object $*$ of the discrete one-object category the arrow f from S to S' . Then $F'(1_*) \circ t(*) = 1_{S'} \circ f = f = f \circ 1_S = t(*) \circ F(1_*)$. Hence t is a natural transformation from F to F' . Conversely, it is easy to show that every natural transformation from F to F' must have this form.

EXAMPLE 34.2

A *small graph* consists of a set of objects, a set of arrows, and two mappings (*source* and *target*) from the set of arrows to the set of objects. A *morphism* F between small graphs is a mapping which sends the objects of the first graph to the objects of the second, and the arrows of the first graph to the arrows of the second. Moreover, if $f : a \rightarrow b$ in the first graph, we require that $F(f) : F(a) \rightarrow F(b)$. That is, graph morphisms preserve source and target. In Example 33.5 we saw that a graph may be viewed as a functor. Now we shall show that a graph morphism can be viewed as a natural transformation between functors which represent graphs.

Let A be the two-object category of Example 33.5 and F any functor from A to the category of sets. We saw that $(Y, X, F(s), F(t))$ forms a small graph, call it G .

Suppose h is a graph morphism from G to $G' = (Y', X', F'(s), F'(t))$, where F' is a second functor from A to the category of sets and $F'(\mathbf{a}) = X'$, $F'(\mathbf{o}) = Y'$. Let τ be a function from the set $\{\mathbf{a}, \mathbf{o}\}$ to the class of arrows in *Set*, which assigns to \mathbf{a} the map $h : X \rightarrow X'$ and to \mathbf{o} the map $h : Y \rightarrow Y'$. To show that τ is a natural transformation, we must show that, if f is either s or t , then $F'(f) \circ \tau(\mathbf{a}) = \tau(\mathbf{o}) \circ F(f)$. (The equations $F'(1_{\mathbf{a}}) \circ \tau(\mathbf{a}) = \tau(\mathbf{a}) \circ F'(1_{\mathbf{a}})$ and $F'(1_{\mathbf{o}}) \circ \tau(\mathbf{o}) = \tau(\mathbf{o}) \circ F'(1_{\mathbf{o}})$ follow at once.)

Both $F'(s) \circ \tau(\mathbf{a})$ and $\tau(\mathbf{o}) \circ F(s)$ map X to Y' . If $a \in X$, then $(F'(s) \circ \tau(\mathbf{a}))(a) = F'(s)(h(a))$, which is the source of $h(a)$ in Y' . Since h is a graph morphism, $F'(s)(h(a)) = h(F(s)(a))$. But this equals $\tau(\mathbf{o}) \circ F(s)(a)$. Similarly, one can show that $F'(t) \circ \tau(\mathbf{a}) = \tau(\mathbf{o}) \circ F(t)$. Thus τ is a natural transformation.

EXAMPLE 34.3

We saw in Chapter 33 that an M -set can also be regarded as a functor. If (M, X, m) and (M, X', m') are two M -sets, an M -homomorphism is a function f from X to X' such that, if $a \in M$ and $b \in X$, $f(m(a, b)) = m'(a, f(b))$. We usually write this equation as $f(ab) = af(b)$. We will let

the reader show that such an M -homomorphism may be viewed as a natural transformation between the functors that represent the M -sets.

It may seem from the above examples that category theory is merely a complicated way of expressing simpler ideas. However, one should bear in mind that the abstract definitions of the theory embody the ideas and methods of many branches of mathematics at once, and thus may serve to unify their separate proofs and results. Chapter 35 should help to illustrate this fact.

Exercises

1. Show that every M -homomorphism can be viewed as a natural transformation.
2. Write an essay supporting one of the following views:
 - (a) Category theory is a perfect example of useless abstraction. Instead of giving us something new in mathematics, it merely burdens us with a new jargon.
 - (b) Category theory is the crown of contemporary mathematics. It combines insights from different branches of mathematics and provides a common language for discussing them.

A Natural Transformation between Vector Spaces

We begin with a review of vector spaces over the field of real numbers, although any other field may be substituted for \mathbf{R} .

A *vector space* V over \mathbf{R} is an Abelian group together with a mapping from $V \times \mathbf{R}$ to V sending (x, r) to xr , such that, if r and s are real numbers and x, y are in V , then

$$(x + y)r = xr + yr,$$

$$x(r + s) = xr + xs,$$

$$x(rs) = (xr)s,$$

$$x1 = x.$$

Note that \mathbf{R} is a vector space over itself.

A *linear transformation* from a vector space V to a vector space V' is a mapping $f : V \rightarrow V'$ such that $f(x + y) = f(x) + f(y)$ and $f(xr) = (f(x))r$ for any $x, y \in V$ and $r \in \mathbf{R}$.

Taking vector spaces as objects and linear transformations as arrows, it is easy to show that the vector spaces (over the reals) form a concrete category which we shall call *Vect*.

A linear transformation with codomain \mathbf{R} is called a *linear functional*. If f and g are linear functionals on V , we define $(f + g)(x) = f(x) + g(x)$ and $(fr)x = f(xr)$. Now the set of linear functionals on V forms a vector space over \mathbf{R} , called the *dual space* V^* of V .

The above procedure may be repeated to obtain the *double dual* of V , namely, $V^{**} = (V^*)^*$. This double dual is closely related to V .

Let $\tilde{} : V \rightarrow V^{**}$ so that, if $x \in V$, then \tilde{x} is the transformation from V^* to \mathbf{R} that maps any linear functional f to $f(x)$, that is, $\tilde{x}(f) = f(x)$. Two things follow immediately:

- I. \tilde{x} is a linear transformation from V^* to \mathbf{R} , that is, a linear functional on V^* ;
- II. $\tilde{}$ is a linear transformation from V to V^{**} .

In the case that V has finite dimension, $\tilde{}$ is an isomorphism.

If h is a linear transformation from a vector space V to a vector space V' , we define h^{**} as the function from V^{**} to V'^{**} such that, if $p \in V^{**}$, $h^{**}(p)$ is the member of V'^{**} that maps f' in V'^* to $p(f' \circ h)$. Note that

III. $f' \circ h : V \rightarrow \mathbf{R}$ and thus $f' \circ h \in V^*$, which is the domain of p ;

IV. h^{**} is a linear transformation from V^{**} to V'^{**} .

In proving IV, we note that, if $f' \in V'^*$, then

$$\begin{aligned} h^{**}(pr)(f') &= (pr)(f' \circ h) = p((f' \circ h)r) \\ &= p((f'r) \circ h) = h^{**}(p)(f'r) = (h^{**}(p)r)(f'). \end{aligned}$$

Now suppose F maps each object V in the category *Vect* to V^{**} , and each arrow $h : V \rightarrow V'$ in *Vect* to h^{**} . If h is the identity function on V , then h^{**} is the identity function on V^{**} , since then $h^{**}(p)(f') = p(f' \circ h) = p(f')$.

Moreover, if $h : V \rightarrow V'$ and $g : V' \rightarrow V''$ are linear transformations, then so is $g \circ h : V \rightarrow V''$. $F(g \circ h) = (g \circ h)^{**}$ maps $p \in V^{**}$ to the member of V''^{**} that maps f'' in V''^* to $p(f'' \circ (g \circ h))$. That is, if $f'' : V'' \rightarrow \mathbf{R}$,

$$(g \circ h)^{**}(p)(f'') = p(f'' \circ (g \circ h)).$$

(Note that $f'' \circ (g \circ h) : V \rightarrow \mathbf{R}$, so that $f'' \circ (g \circ h) \in V^*$, which is the domain of $p \in V^{**}$.)

Since

$$\begin{aligned} (F(g) \circ F(h))(p)(f'') &= g^{**}(h^{**}(p))(f'') = h^{**}(p)(f'' \circ g) \\ &= p((f'' \circ g) \circ h) = p(f'' \circ (g \circ h)), \end{aligned}$$

it follows that F is a functor from *Vect* to *Vect*.

Another functor from *Vect* to *Vect* is the identity functor I .

Let t assign to every vector space V the linear transformation from V to $F(V) = V^{**}$ which we called $\tilde{}$. That is, let $t(V)(x) = \tilde{x}$. Suppose $h : V \rightarrow V'$ and let x be any element of V . Then $(F(h) \circ t(V))(x) = h^{**}(\tilde{x})$.

Also, $(t(V') \circ I(h))(x) = (h(x))^\sim$. These two elements of V'^{**} are in fact equal. For let $f' : V' \rightarrow \mathbf{R}$, so that $f' \in V'^*$. Then

$$\begin{aligned} h^{**}(\tilde{x})(f') &= \tilde{x}(f' \circ h) \\ &= (f' \circ h)(x) \\ &= f'(h(x)) \\ &= (h(x))^\sim(f'). \end{aligned}$$

We may conclude that $F(h) \circ t(V) = t(V') \circ I(h)$, and hence t is a natural transformation from the functor I to the functor F .

Examples such as this have led to the slogan *that many objects of interest in mathematics are functors and that the arrows between them are natural transformations*. This and the slogans mentioned earlier were first proposed by F. W. Lawvere.

Exercises

1. Show that *Vect* is a concrete category.
2. Show that the sum of two linear transformations (from V to V') is a linear transformation.
3. Show that V^* is a vector space.
4. Verify **I**, **II**, **III** and **IV** from the text.
5. Generalize the results of this chapter from vector spaces to M -sets. (Things become a little easier if it is assumed that multiplication in M is commutative.)

References

- [1] *Dictionary of Scientific Biography*. New York: Charles Scribner's Sons, 1976.
- [2] Allen, Reginald E. *Greek Philosophy: Thales to Aristotle*. New York: The Free Press, 1966.
- [3] Altmann, Simon L. 'Hamilton, Rodrigues, and the Quaternion Scandal'. *Mathematics Magazine*, 62 (1989) 291-308.
- [4] Anbona, S. 'Un Traité d'Abu Jafar [al Kazin] sur les Triangles Rectangles Numériques'. *J. History Arabic Studies* 3 (1979), 134-178.
- [5] Anglin, W. S. *Mathematics: A Concise History and Philosophy*, New York: Springer-Verlag, 1994.
- [6] Anglin, W. S. *The Queen of Mathematics*, Dordrecht: Kluwer, 1995.
- [7] Anglin, W. S. 'Using Pythagorean Triangles to Approximate Angles'. *American Mathematical Monthly*, 95 (1988), 540-41.
- [8] Aquinas, Thomas *Summa Contra Gentiles*. Trans. Anton C. Pegis et al. London: University of Notre Dame Press, 1975.
- [9] Archibald, R. C. 'Gauss and the Regular Polygon of Seventeen Sides'. *American Mathematical Monthly*, 27 (1920), 323-26.
- [10] Archimedes *The Works of Archimedes*. Trans. T. L. Heath. New York: Dover, 1897.

- [11] Aryabhata *Aryabhatiya of Aryabhata*, New Delhi: Indian National Science Academy, 1976.
- [12] Aschenbrenner, Karl *The Concept of Value*. Dordrecht: D. Reidel, 1971.
- [13] Ascher, M. *Ethnomathematics*. Pacific Grove, California: Brooks/Cole, 1991.
- [14] Ayoub, R. 'What is a Napierian Logarithm?'. Amer. Math. Monthly, 100 (1993), 351-364.
- [15] L. Bakhtiar *Sufi Expressions of the Mystic Quest*. London: Thamer and Hudson, 1976.
- [16] Ball, W. W. Rouse *A Short Account of the History of Mathematics*. New York: Dover, 1960.
- [17] Barker, Stephen F. *Philosophy of Mathematics*. Englewood Cliffs: Prentice-Hall, 1964.
- [18] Bailey, Cyril *The Greek Atomists and Epicurus*. Oxford: Clarendon, 1928.
- [19] Barnes, Jonathan *Early Greek Philosophy*. London: Penguin, 1987.
- [20] Barnes, Jonathan *PreSocratic Philosophers*. Vol. 2. London: Routledge and Paul, 1979.
- [21] Barnsley, Michael *Fractals Everywhere*. Boston: Academic, 1988.
- [22] Bell, E. T. *The Development of Mathematics*. 2nd ed. New York: McGraw-Hill, 1945.
- [23] Bell, E. T. *Men of Mathematics*. New York: Simon and Schuster, 1965.
- [24] Benacerraf, Paul, and Hilary Putnam, ed. *Philosophy of Mathematics*. 2nd ed. Cambridge: Cambridge University, 1983.
- [25] Berggren, J. L. *Episodes in the Mathematics of Medieval Islam*. New York: Springer-Verlag, 1986.
- [26] Birkhoff, Garret, and Saunders Mac Lane *A Survey of Modern Algebra*. 4th ed. New York: Macmillan, 1977.
- [27] Bishop, Errett, and Douglas Bridges *Constructive Analysis*. Berlin: Springer-Verlag, 1985.
- [28] Bleicher, M. N. 'A New Algorithm for the Expansion of Egyptian Fractions'. Journal of Number Theory, 4 (1972), 342-82.

- [29] Blumenthal, L. M. *A Modern View of Geometry*. San Francisco: Freeman, 1961.
- [30] Boyer, Carl B. and Uta C. Merzbach *A History of Mathematics*. 2nd ed. New York: John Wiley, 1989.
- [31] Browder, F. E. (ed.) *Mathematical Developments Arising from Hilbert Problems*. Proc. Symposia Pure Math. 28, Providence: Amer. Math. Soc., 1976.
- [32] Bruckheimer, M. and Y. Salomon 'Some Comments on R. J. Gillings' Analysis of the $2/n$ Table in the Rhind Papyrus'. *Historia Mathematica*, 4 (1977), 445-52.
- [33] Burton, David M. *The History of Mathematics*. Dubuque: Wm. C. Brown, 1985.
- [34] Cajori, Florian *A History of Mathematics*. 2nd ed. New York: Macmillan, 1919.
- [35] Cantor, Georg *Transfinite Numbers*. New York: Dover, 1955
- [36] Cardano, Girolamo *The Great Art*. Trans. T. Richard Witmer Cambridge, Mass.: MIT, 1968.
- [37] Cardano, Girolamo *The Book of My Life*. Trans. Jean Stoner, New York: Dover, 1962.
- [38] Carruccio, Ettore *Mathematics and Logic in History and in Contemporary Thought*. Trans. I. Quigly. London: Faber and Faber, 1964.
- [39] Carslaw, H. S. 'Gauss's Theorem on the Regular Polygons which can be Constructed by Euclid's Method'. *Proceedings of the Edinburgh Mathematical Society*, 28 (1910), 121-8.
- [40] Chace, Arnold Buffum et al. *The Rhind Mathematical Papyrus*. Oberlin, Ohio: Mathematical Association of America, 1927-9.
- [41] F. Chung and S. Sternberg 'Mathematics and the Buckyball'. *Amer. Scientist* 81 (1993).
- [42] Church, A. *The Calculi of Lambda Conversion*. Ann. Math. Studies 6, Princeton University, 1941.
- [43] Connell, I. *Modern Algebra*. New York: North Holland, 1982.
- [44] Couture, J. and J. Lambek 'Philosophical Reflections on the Foundations of Mathematics'. *Erkenntniss* 34 (1991), 187- 209.
- [45] Coxeter, H. S. M. and S. L. Greitzer *Geometry Revisited*. Washington: The Mathematical Association of America, 1967.

- [46] Dauben, J. W. 'Abraham Robinson and Nonstandard Analysis'. in: W. Asprey and P. Kitcher (eds.), *Studies in the Philosophy of Science* 11 (1988), Minneapolis: University Minneapolis, 1988.
- [47] Dauben, J. W. *Georg Cantor: His Mathematics and Philosophy*. Cambridge, Mass: Harvard University, 1979.
- [48] Davis, Martin 'Hilbert's Tenth Problem is Unsolvble'. *American Mathematical Monthly*, 80 (1973), 233 -69.
- [49] Descartes, René *The Geometry of Rene Descartes*. Trans. D. E. Smith and M. L. Latham. New York: Dover, 1954.
- [50] L. E. Dickson *History of the Theory of Numbers*. New York: Chelsea, 1971.
- [51] Diophantus *Arithmetica*. 2nd ed. Trans. Thomas L. Heath. New York: Dover, 1964.
- [52] Diophantus *Books IV to VII of Diophantus' Arithmetica*. Trans. Jacques Sesiano. New York: Springer-Verlag, 1982.
- [53] Dodgson, Charles L. *Euclid and his Modern Rivals*. London: Macmillan, 1879.
- [54] Dummett, M. *Elements of Intuitionism*. Oxford: Oxford University Press, 1977.
- [55] Dunham, William *Journey through Genius*. New York: Wiley, 1990.
- [56] Eilenberg, S. and S. Mac Lane 'General Theory of Natural Equivalence'. Trans. Amer. Math. Soc., 58 (1945) 231-94.
- [57] Ellison, W. J. 'Waring's Problem'. Amer. Math. Monthly 73 (1971), 10-36.
- [58] Euclid *Elements*. Trans. T. L. Heath. 2nd ed. New York: Dover, 1956.
- [59] Euler, Leonhard *Opera Omnia*. Geneva: 1944.
- [60] Eves, Howard *An Introduction to the History of Mathematics*. 5th ed. New York: Holt, Rinehart and Winston, 1969.
- [61] Eves, Howard *A Survey of Geometry*. Boston: Allyn and Bacon, 1963.
- [62] Ewald, Gunter *Geometry: An Introduction*. Belmont, California: Wadsworth, 1971.
- [63] Federico, P. J. *Descartes on Polyhedra*. New York: Springer-Verlag, 1982.

- [64] Fermat, Pierre *Oeuvres de Fermat*. Paris: Gauthier-Villar, 1891.
- [65] Fibonacci *Le Livre des Nombres Carres*. Trans. Paul ver Eecke. Paris: Albert Blanchard, 1952.
- [66] Fowler, D. H. *The Mathematics of Plato's Academy*. Oxford: Clarendon, 1987.
- [67] Freed, S. and R. S. Freed 'Origin of the Swastika', *Natural History*, 1981, 68-75.
- [68] Gardner, Martin 'Mathematical Games' *Scientific American*, 233 (December 1975), 117-8.
- [69] Gillings, Richard J. *Mathematics in the Time of the Pharaohs*. New York: Dover, 1972.
- [70] Gittleman, Arthur *History of Mathematics*. Columbus: Charles E. Merrill, 1975.
- [71] Glass, A. M. W. 'Existence Theorems in Mathematics'. *The Mathematical Intelligencer*, 11 (1989), 56-62.
- [72] Goodstein, R. L. *Recursive Number Theory*. New York: Academic, 1970.
- [73] Graves, Robert Perceval *Life of Sir William Rowan Hamilton*. Dublin: Hodges, Friggs Co., 1885.
- [74] Guthrie, Kenneth Sylvan *The Pythagorean Sourcebook and Library*. Grand Rapids: Phanes, 1987.
- [75] Hall, H. S. and Knight, S. R. *Higher Algebra*. London: Macmillan, 1955.
- [76] Hardy, G. H. *A Mathematician's Apology*. Cambridge: Cambridge University, 1992.
- [77] Hatcher, William S. *The Logical Foundations of Mathematics*. Oxford: Pergamon, 1982.
- [78] Hawkins, G. S. *Stonehenge Decoded*. New York: Doubleday, 1965.
- [79] Heaslett, M. A. and J. V. Uspensky *Elementary Number Theory*. New York: McGraw Hill, 1939.
- [80] Heath, T. L. *A History of Greek Mathematics*. Oxford: Clarendon, 1921.
- [81] Heuer, K. *City of Stargazers*. New York: Charles Scribner's Sons, 1972.

- [82] Hilbert, D. *The Foundations of Geometry*. La Salle: Open Court, 1962.
- [83] Hintikka, J. *The Philosophy of Mathematics*. Oxford: Oxford University, 1969.
- [84] Hofstadter, D. R. *Gödel, Escher, Bach*. New York: Basic Books, 1979.
- [85] Hooper, Alfred *Makers of Mathematics*. New York: Random House, 1948.
- [86] Hoyle, F. *On Stonehenge*. San Francisco: Freeman, 1977.
- [87] Huntley, H. E. *The Divine Proportion*. New York: Dover, 1970.
- [88] Jones, P. James 'Diophantine Representations of Fibonacci Numbers over Natural Numbers'. In *Applications of Fibonacci Numbers*. Vol. 3. Ed. G. E. Bergum et al. Dordrecht: Kluwer, 1990.
- [89] Jones, Philip S. 'Recent Discoveries in Babylonian Mathematics 1'. *The Mathematics Teacher*, 50 (1957), 162-5.
- [90] Johnson, Roger A. *Advanced Euclidean Geometry*. New York: Dover, 1960.
- [91] Jung, Karl *Man and His Symbols*. New York: Doubleday, 1983.
- [92] Kershner, R. B., and L. R. Wilson *The Anatomy of Mathematics*. New York: Ronald, 1950.
- [93] Kershner, R. B. 'On Paving the Plane'. *American Mathematical Monthly*, 75 (1968), 839-44.
- [94] Khayyam, O. *The Algebra of Omar Khayyam*. Ed. D. S. Kasir New York: AMS, 1972.
- [95] Khayyam, O. *Rubaiyat*. New York: Doubleday, 1952.
- [96] Khwarizmi, M. *The Algebra*. Trans. Frederic Rosen. London: Oriental Translation Fund, 1831.
- [97] Klarner, David A., ed. *The Mathematical Gardner*. Boston: Prindle, Weber and Schmidt, 1980.
- [98] Kleene, S. C. *Intoduction to Metamathematics*. New York: Van Nostrand, 1952.
- [99] Kline, Morris *Mathematics in Western Culture*. New York: Oxford University, 1953.
- [100] Kneale, W. C., and M. Kneale *The Development of Logic*. Oxford: Clarendon, 1984.

- [101] Koblitz, Neal *A Course in Number Theory and Cryptography*. New York: Springer-Verlag, 1987.
- [102] Koestler, Arthur *The Watershed*. New York: Anchor Books, Doubleday, 1960.
- [103] Lagrange, J. *Oeuvres de Lagrange*. Paris: Gauthier-Villars, 1869.
- [104] Lam, C. W. H. 'How Reliable is a Computer-Based Proof?' *The Mathematical Intelligencer*, 12(1) (1990), 8-12.
- [105] Lambek, J. 'Are the Traditional Philosophies of Mathematics Really Incompatible?'. *The Mathematical Intelligencer*, 16 (1994), 56-62.
- [106] Lambek, J. 'If Hamilton had Prevailed: Quaternions in Physics'. *The Mathematical Intelligencer*, to appear.
- [107] Lambek, J. 'How to Program an Abacus?'. *Can. Math. Bulletin*, 4 (1961), 295-302.
- [108] Lambek, J. and P.J. Scott *Introduction to Higher Order Categorical Logic*. Cambridge: Cambridge University, 1988.
- [109] Legendre, Adrien-Marie *Theorie des Nombres*. 4th ed. Paris: Firmin Didot, 1830.
- [110] LeLionnais, Francois, ed. *Great Currents of Mathematical Thought*. Trans. Charles Pinter and Helen Kline, New York: Dover, 1971.
- [111] Leonard de Pise. See Fibonacci.
- [112] Lipschitz, M. 'Recherches sur la transformation...'. *Journal de Mathematiques Pures et Appliquees* 4th Ser. Vol. 2 (1886) 373-439, especially p. 404.
- [113] Mac Lane, S. *Categories for the Working Mathematician*. New York: Springer-Verlag, 1971.
- [114] Marius, Cleyet-Michaud *Le Nombre d'Or*. Que-Sais-je?, 1530.
- [115] Martin, Georges E. *Transformation Geometry*. New York: Springer-Verlag, 1982.
- [116] Maxwell, J. C. 'Remarks on the Classification of Physical Quantities'. *Proc. London Math. Soc.* 3 (1869), 224-232.
- [117] Maziarz, Edward A. and Thomas Greenwood *Greek Mathematical Philosophy*. New York: Frederick Ungar, 1968.
- [118] McClenon, R. B. 'Leonardo of Pisa and his Liber Quadratorum'. *American Mathematical Monthly*, 26 (1919), 1-8.

- [119] Menninger, K. *Number Words and Number Systems; a cultural history of numbers*. Cambridge, Mass.: M.I.T., 1969.
- [120] Midonick, Henrietta O., ed. *The Treasury of Mathematics*. New York: Philosophical Library, 1965.
- [121] Mohanty, S. P. 'Integer Points of $y^2 = x^3 - 4x + 1$ '. *Journal of Number Theory*, 30 (1988), 86-93.
- [122] Mordell, L. J. *Diophantine Equations*. New York: Academic, 1969.
- [123] Nagel, E. and J.R. Newman *Gödel's Proof*. London: Routledge and Kegan Paul, 1959.
- [124] Nahin, P. J. 'Oliver Heaviside'. *Scientific American*, 1990, 122-129.
- [125] Napier, John *The Construction of the Wonderful Canon of Logarithms*. Trans. William Rae MacDonald. Edinburgh: William Blackwood and Sons, 1889.
- [126] Nelson, Harry L. 'A Solution to Archimedes' Cattle Problem'. *Journal of Recreational Mathematics*, 13 (1980-81), 164-76.
- [127] Neugebauer, O. *The Exact Sciences in Antiquity*. 2nd ed. New York: Dover, 1969.
- [128] Newman, James R., ed. *The World of Mathematics*. New York: Simon and Schuster, 1956.
- [129] Ogilvy, C. Stanley *Excursions in Geometry*. New York: Oxford University, 1969.
- [130] Ore, Oystein *Cardano, the Gambling Scholar*. New York: Dover, 1965.
- [131] Ore, Oystein *Number Theory and Its History*. New York: McGraw-Hill, 1948.
- [132] Pappus of Alexandria *La Collection Mathématique*. Trans. Paul ver Eecke, Paris: Bruges, 1933.
- [133] Peano, Ioseph *Arithmetices Principia*. Rome: Fratres Bocca, 1889.
- [134] Peitgen, H. O. and P. H. Richter *The Beauty of Fractals*. Berlin: Springer-Verlag, 1986.
- [135] Penrose, Roger *The Emperor's New Mind*. New York: Oxford University, 1989.
- [136] Plato *The Collected Dialogues*. New Jersey: Princeton University, 1989.

- [137] Plutarch, *Makers of Rome*. London: Penguin, 1965.
- [138] Poincaré, H. *Science and Hypothesis*. New York: Walter Scott, 1905.
- [139] Prawitz, D. *Natural Deduction*. Stockholm: Almquist and Wiskell, 1965.
- [140] Robins, Gay, and Charles Shute, ed. *The Rhind Mathematical Papyrus*. London: British Museum Publications, 1987.
- [141] Rosen, F. *The Algebra of Mohammed Ben Musa*. London: 1831.
- [142] Rosenbloom *The Elements of Mathematical Logic*. New York: Dover, 1950.
- [143] Russell, Bertrand *A History of Western Philosophy*. New York: Simon and Schuster, 1972.
- [144] Russell, Bertrand *Mysticism and Logic*. London: Penguin Books, 1953.
- [145] Shapiro, Harold N. *Introduction to the Theory of Numbers*. New York: John Wiley, 1983.
- [146] Sierpinski, W. *Elementary Theory of Numbers*. Warsaw: Polska Akademia Nauk, 1964.
- [147] Silberstein, L. *The Theory of Relativity*. 2nd ed. London: Macmillan, 1924.
- [148] Smith, David Eugene *History of Mathematics*. New York: Dover, 1958.
- [149] Smith, David Eugene *A Source Book in Mathematics*. New York: Dover, 1959.
- [150] Smith, T. V., ed. *From Thales to Plato*. Chicago: The University of Chicago, 1965.
- [151] Stillwell, John *Mathematics and Its History*. New York: Springer-Verlag, 1989.
- [152] Tymoczko, Thomas *New Directions in the Philosophy of Mathematics*. Boston: Birkhauser, 1985.
- [153] Tzanakis, N. and B. M. M. de Weger 'On the Practical Solution of the Thue Equation'. *Journal of Number Theory*, 31 (1989), 99-132.
- [154] J.V. Uspensky and M.A. Heaslet *Elementary Number Theory*. McGraw Hill, New York: 1939.

- [155] Uspensky, V. A. *Post's Machines*. Moscow: Mir Publications, 1979.
- [156] Van der Waerden, B. L. *Geometry and Algebra in Ancient Civilizations*. Berlin: Springer-Verlag, 1983.
- [157] Van der Waerden, B. L. *A History of Algebra*. Berlin: Springer-Verlag, 1985.
- [158] Van der Waerden, B. L. *Science Awakening*. Vol. 1. 4th ed. Trans. Arnold Dresden, Dordrecht: Kluwer, 1975.
- [159] van Nooten, B. 'Binary Numbers in Indian Antiquity'. *J. Indian Philosophy* 21 (1993), 31-50.
- [160] Ver Eeke, P. *Diophante d'Alexandrie*. Paris: Blanchard, 1959.
- [161] Viète, Francois *Opera Mathematica*. Hildesheim: Georg Olms Verlag, 1970.
- [162] Wantzel, P. L. 'Recherches sur les moyens de reconnaitre si un probleme de geometrie peut se resoudre par la regle et le compas'. *Journal de Mathematique*, 2 (1837) 366-72.
- [163] Weyl, André *Number Theory: An Approach through History*. Basel: Birkhauser, 1984.
- [164] Weyl, Herman *Space, Time and Matter*. New York: Dover, 1922.
- [165] Zaslavski, C. *Africa Counts*. Boston: Prindle, Weber and Schmidt, 1973.

Index

- abacus 247
- Abel 136
- Abelian 184
- Abraham 21
- abstraction 259, 295, 305
- Abu'l Wafa 165
- Academy 63, 109
- Achilles 55
- Adelhard 121
- Agha Khan 119
- Ahmose 7, 11
- aleph-null 68, 219
- Alexander the Great 83
- Alexandria 29, 83, 93, 95, 107
- algebra 106, 118
- algorithm 118, 245, 256
- al-Khayyami 118
- al-Khazin 106, 122, 196
- al-Khwarizmi 117
- al-Mamun 117
- amicable 34, 119
- Amru 109
- Anaximander 31
- Anaximenes 31
- Antiphon 60
- Apollonius 94, 130
- Aquinas 84, 109
- arbelos 101
- Archimedes 27, 97, 104, 112, 141, 153
- Archytas 63, 84
- Aristarchus 93
- Aristotle 7, 33, 47, 53, 68, 84, 99, 100, 156, 265, 266, 279
- Arius 107
- arrow 55, 162, 297
- Aryabhata 113
- Asia Minor 2
- astrology 105, 114, 129, 130, 143
- astronomy 1, 105, 141
- Aswan 95
- Athanasius 107
- Athens 59
- atom 57, 69
- Augustine 39
- axiom 84, 289
- axiom of Archimedes 92, 99, 146
- axiom of choice 285
- Babylonian 21

- Bachet 149
 backwards causation 54
 Baghdad 115
 Barrow 155
 Beltrami 90
 Ben Ezra 121
 Berkeley 56, 160
 Bernays 92, 268, 301
 Bhaskara 113, 114
 Bible 109, 126
 binary 11
 Binomial Theorem 155
 biquaternion 212
 Boethius 109
 Bolyai 90
 Bombelli 130
 Book of Changes 11
 Boole 267
 Bourbaki 83, 84, 87, 175
 Brahe 142
 Brahmagupta 3, 113, 233
 Briggs 141
 Brouwer 176, 271, 277
 buckminster fullerene 45
 buckyball 45

 calculation 245
 calculus 98, 118, 156, 281
 Cantor 68, 219
 Cardano 3, 129, 195
 cardinal 219
 cartesian 90, 170
 Casanova 80
 category 162, 295, 297
 cattle problem 100
 Cauchy 191
 Cavalieri 153
 chemical equation 42
 Chen Jing-Run 18, 112
 China 3, 11, 18, 111, 125, 149, 151
 Ch'in Chiu Shao 112
 Chiu Chang Suan Shu 111
 Chinese remainder 111
 chord 104, 113, 130, 153
 Christian 112
 Christina 146
 Ch'ung Chih 112
 Church 147, 176, 252, 259
 Chu Shih Chieh 112
 cipher key 18
 circle 9, 60, 64, 98, 153
 circle squaring 60, 71, 79
 classical problems 71, 79
 Cleopatra 83
 Cohen 220, 245
 Cole 38
 Columbus 95, 126
 combinator 262
 commensurable 47
 complete ordered field 192
 completeness 86, 289, 290
 complex 3, 122, 129, 130, 176, 195
 comprehension scheme 267
 concrete category 295
 conic section 64, 94, 147
 Constantine 107
 construction, ruler and compass 71
 continuous 30, 61, 85, 199
 continued fractions 228
 continuum hypothesis 245
 convergent 165, 228
 Copernicus 1, 94, 141
 cosines, law of 86
 count nouns 30
 Couture 285
 Coxeter 151, 165
 creator 39, 43
 Crotone 33
 cryptography 18
 cubic equation 118, 119, 128, 132
 Curry 224
 cycloid 154

 da Coi 128
 d'Alembert 199
 Davis 246
 Dedekind 3, 269

- De Gang Ma 112
- de la Vallée Poussin 18
- Delian problem 62, 71
- Democritus 7, 31, 57
- De Moivre 73, 124, 164, 165, 197
- De Morgan 267, 286
- Desargues 145
- Descartes 44, 75, 146, 164
- Dieudonné 87, 175, 176
- Diophantine 100, 233, 255
- Diophantus 27, 100, 105, 149, 233
- Dirac 216
- directrix 95
- Dirichlet 149
- discrete 30, 61
- distributive law 48
- divine 34, 126
- division ring 204
- double negative 265, 271, 275
- doubling cube 62, 65, 71, 80, 81, 147
- eccentricity 95
- eclipse 2, 29
- Egypt 7, 16, 41, 159
- Eilenberg 303
- Einstein 212, 217
- electromagnetic 211, 217
- electron 54
- Elements* 16, 61, 83, 121
- ellipse 153
- elliptic geometry 85, 91
- Empedocles 57
- enumerable 251
- Eratosthenes 17, 95
- Erdős 10, 16, 245
- Euclid 16, 34, 37, 48, 83, 156, 196, 228
- Eudoxus 3, 61, 63–64, 84, 86, 191
- Euler 39, 44, 82, 164, 207
- Euler's Formula 44, 165
- excluded middle 68, 100, 265, 271
- Ezekiel 22
- factorization 16, 18
- Faltings 150
- Fermat 81, 148, 165, 242
- Ferrari 129
- Ferro 128
- Feynman 54
- Fibonacci 106, 122, 256
- field 80, 87, 162, 184, 189
- figurative 35, 169
- Fiore 128
- Fitzgerald 119
- fixpoint 262
- flow diagram 247
- fluent 157
- fluxion 155, 163
- focus 95
- formalist 176, 294
- Four Colour Theorem 295
- Four Weight Problem 12
- Fowler 227
- Fraenkel 268
- Francis 122
- Frederick II 122, 166
- Frege 267
- Frey 243
- function 259, 262, 296
- functor 299
- Fundamental Theorem of Algebra 199
- Fundamental Theorem of Arithmetic 16, 231
- Galileo 141, 147, 156, 320
- Galois 136
- Gauss 3, 16, 81, 86, 90, 169, 199
- generating function 123
- Gentzen 273
- George I 160
- Gerbert 121
- Girard 199
- God 22, 148, 150, 151, 161, 176, 268
- Gödel 220, 245, 253, 268, 289
- Goldbach 17, 18, 112

golden ratio 64, 74, 124, 126,
230, 237

Golenishchev 7

graha 1

graph 297, 304

gravity 156

Gregory 154, 155

group 162, 184

Hadamard 18

Hafiz 119

Hamilton 3, 175, 203

Harriot 126

Hasan 119

Heath 106

Heaviside 211

Hegel 53

Henkin 289

Heraclitus 31, 53

Hermite 79, 212

Herodotus 7

Heron 22, 104, 114

Heyting 277

Hilbert 17, 87, 91, 107, 146, 245,
255, 293

Hipparchus 22, 105

Hippasus 43, 47

Hippias 59

Hippocrates 61, 84

Hiram of Tyre 27

Hooke 156

Horner 112

Hoyle 2

Hrotsvitha 121

Huygens 154

Hypatia 108

hyperbolic geometry 90

icosahedron 42, 44

India 3, 112, 155, 233

infinite 31, 68, 114, 118, 219

infinitesimal 56, 160

injective 189

integers 183

integral domain 183

intuitionist 68, 176, 272, 294

intuitionistic logic 271, 281

Ionia 29, 69

Iraq 2, 21

irrationality 35, 47

Jesus 107, 129, 139

Jones 98

Jordanus 122

Joseph 8

Jupiter 1, 142

Justinian 109

kebu 1

Kepler 142, 156

Khayyam 119, 128

Kingsley 108

Kleene 251

Klein-Gordon 215

Koerberero 27

Kolmogorov 277

Kronecker 28, 176, 268

Kummer 149

Lagrange 113, 166, 207, 237, 255

Lambda Calculus 259, 289

Laplace 167

Lawvere 54, 309

Least Action 161

Legendre 17, 86, 89, 149, 167,
170

Leibniz 56, 159, 267

Lenin 54

Leonardo, see Fibonacci

Lewis 109

Lilavati 114

limit 55, 160

Lindemann 79

Lipschitz 208

linear equations 111, 114, 233

linear functional 307

linear transformation 307

Lobachevsky 90, 203

logarithm 17, 139

logic 53, 265

- logician 176, 267, 294
- Lorentz transformation 212
- Louis XIV 159
- love 57
- Lucas 38
- lune 61
- Mac Lane 303
- Maclaurin 163
- magic square 111
- Magna Grecia 2
- Mahavira 114
- Marcellus 98
- Marie Antoinette 166
- Mars 1, 142, 143
- Marx 53, 57
- mass nouns 30
- mathematical induction 61, 151, 179
- Matiyasevič 17, 123, 246, 257
- matrix 203
- Maupertuis 161
- Maxwell 211
- Mazur 243
- means 34
- Menaechmus 64
- Menelaus 104, 107
- Meno 63
- Mercator 154
- Mercury 1
- Mersenne 38, 145, 148
- mesopotamian 21, 122, 242
- M-homomorphism 304
- Miletus 29
- Mill 161
- Minkowski 212
- model 289
- monad 162
- monoid 296, 298, 300
- month 1
- moon 1, 104, 166
- Moscow Papyrus 7
- M-set 300
- Museum 83
- mysticism 39, 162
- Napier 139
- natural number 175, 179
- natural transformation 303, 307
- Nebuchadnezzar 21
- negative number 3, 113, 123, 126
- Nehemiah 28
- Nelson 101
- Newton 56, 84, 97, 148, 155
- Nicomachus 108
- Nim 12
- Nizam-ul-Mulk 119
- nodes 1, 297
- Noether 295
- nonconstructive 271, 272
- non-Euclidean 89
- norm 122, 204
- notation 125
- Omar I 109
- omega-complete 290
- omega-consistent 293
- Oresme 122
- Pacioli 126
- Paganini 39
- Pappus 107, 146
- papyrus 7
- paradox 55, 268
- Parallel Postulate 85, 89, 118, 203
- Parmenides 54
- Pascal 112, 150
- Pauli 215
- Peano 176, 179, 182
- Pell 113, 160, 237
- pentagon 42
- perfect 37, 86, 121, 122, 145, 164
- Philo of Alexandria 39
- Philo of Megara 266
- pi 9, 27, 77, 98, 112, 113, 114, 118, 155
- Pingala 112
- Planck 215
- planet 1

- Plato 29, 42, 50, 59, 60, 63, 67, 70, 162
 Platonist 68, 176, 290, 294
 Playfair 89
 Plimpton 27, 33
 Plotinus 105
 Plutarch 98
 Poincaré 90, 213, 269
 polyhedron, *see* regular solid
 positron 54
 Post 246
 postulate 84
 prime 15, 18, 81, 86, 257
 primitive recursive 246
 primitive root 169
 probability 129, 151, 164
 Proclus 50, 89, 109
 program 247
 projective geometry 145
 proof 30
 Ptolemy 22, 105, 108, 130, 141
 Ptolemy I 83
 Putnam 246, 257
 pyramid 7, 31, 57, 113
 Pythagoras 27, 31, 33, 43
 Pythagoreans 3, 33, 47, 108, 118, 142, 237
 Pythagorean triangle 27, 241

 Qin Jiushao 112
 quadratic equation 25
 quadratic reciprocity 168, 169
 quadratic surds 237
 quadratrix 60
 quadrivium 34
 quantum mechanics 155, 215
 quantum of time 56
 quartic equation 129, 135
 quaternions 3, 175, 203

 rahu 1
 Ramanujan 115
 rational operation 75
 rational 187
 real 3, 191

 Recorde 126
 recursive 251
 Regiomontanus 125
 regular polygon 41, 60, 73, 81, 86
 regular solid 41, 63, 87, 142
 relativity 212
 Renaissance 125
 Rhind Papyrus 7
 Ribet 243
 Riemann 91
 Riese 126
 rigour 86, 87
 ring 183
 Robinson, A. 56, 160, 162
 Robinson, J. 246, 255
 Rodrigues 3
 Rudolff 126
 Ruffini 136
 rule of signs 147
 Russell 162, 176, 267
 Rutherford 57

 Sacchieri 89
 Saturn 1, 142
 scale (base) 9, 11
 Schanuel 192
 Schönfinkel 262
 Schrödinger 215
 Schroeder-Berstein 225
 Schwenter 227
 Shimura 243
 Shyreswood 266
 siddhanta 113
 Silberstein 212
 sine 104, 113
 Six Weight Problem 12
 skew field 204
 Snell 161
 Socrates 59, 62, 162
 spherical triangle 126
 Spinoza 84
 square 47, 63, 67, 86, 207
 squaring circle, *see* circle
 squaring
 Steiner 91

Stevin 141
 Stifel 126, 139
 Stonehenge 2
 string 47, 176, 217, 252, 289, 292
 subroutine 248
 Sultana 235
Sulvasutras 112
 Sumeria 21
 Sun 1, 93
 Sun Tsu 111
 Syene 95
 syllogism 266
 Sylvester 121
 Syracuse 97

 Taniyama 243
 Tartagliā 128
 Taurinus 90
 Taylor, B. 163
 Taylor, R. 243
 Thabit 39, 118
 Thales 3, 29, 30, 33
 Theaetetus 63, 84, 237
 Theodorus 63
 Thibault 89
 Timaeus 42, 63
 tortoise 55
 Tournesol 177, 241, 243
 trigonometry 60, 104, 118
 trisection of angle 60, 71, 130
 Trotsky 47
 Turing 246, 259
 type 268, 285
 Tzanakis 36

unique factorisation 16, 169, 231
 unit fraction 9, 100, 230
 Ur 21

 vector 211, 307
 Venus 1
 Viète 130
 Vinogradov 18
 Voltaire 156, 161
 Von Neumann 176

 Wallis 89, 148, 154, 155
 Wantzel 60, 65
 water 30, 31, 42, 57
 week 1
 Weierstrass 160, 199
 Whitehead 268
 Widman 126
 Wiles 149, 243
 William of Shyreswood 266
 Williams 101
 Wilson 166
 Wittgenstein 266
 Wren 156

 year 1

 Zeno 54, 61
 Zermelo 268
 zero 105, 112, 114, 121, 185
 Zhu Shijie 112
 ziggurat 27

Undergraduate Texts in Mathematics

Anglin: Mathematics: A Concise History and Philosophy.

Readings in Mathematics.

Anglin/Lambek: The Heritage of Thales.

Readings in Mathematics.

Apostol: Introduction to Analytic Number Theory. Second edition.

Armstrong: Basic Topology.

Armstrong: Groups and Symmetry.

Bak/Newman: Complex Analysis.

Banchoff/Wermer: Linear Algebra Through Geometry. Second edition.

Berberian: A First Course in Real Analysis.

Brémaud: An Introduction to Probabilistic Modeling.

Bressoud: Factorization and Primality Testing.

Bressoud: Second Year Calculus.
Readings in Mathematics.

Brickman: Mathematical Introduction to Linear Programming and Game Theory.

Cederberg: A Course in Modern Geometries.

Childs: A Concrete Introduction to Higher Algebra. Second edition.

Chung: Elementary Probability Theory with Stochastic Processes. Third edition.

Cox/Little/O'Shea: Ideals, Varieties, and Algorithms.

Croom: Basic Concepts of Algebraic Topology.

Curtis: Linear Algebra: An Introductory Approach. Fourth edition.

Devlin: The Joy of Sets: Fundamentals of Contemporary Set Theory. Second edition.

Dixmier: General Topology.

Driver: Why Math?

Ebbinghaus/Flum/Thomas: Mathematical Logic. Second edition.

Edgar: Measure, Topology, and Fractal Geometry.

Fischer: Intermediate Real Analysis.

Flanigan/Kazdan: Calculus Two: Linear and Nonlinear Functions. Second edition.

Fleming: Functions of Several Variables. Second edition.

Foulds: Combinatorial Optimization for Undergraduates.

Foulds: Optimization Techniques: An Introduction.

Franklin: Methods of Mathematical Economics.

Hairer/Wanner: Analysis by Its History.
Readings in Mathematics.

Halmos: Finite-Dimensional Vector Spaces. Second edition.

Halmos: Naive Set Theory.

Hämmerlin/Hoffmann: Numerical Mathematics.

Readings in Mathematics.

Iooss/Joseph: Elementary Stability and Bifurcation Theory. Second edition.

Isaac: The Pleasures of Probability.
Readings in Mathematics.

James: Topological and Uniform Spaces.

Jänich: Linear Algebra.

Jänich: Topology.

Kemeny/Snell: Finite Markov Chains.

Kinsey: Topology of Surfaces.

Klambauer: Aspects of Calculus.

Lang: A First Course in Calculus. Fifth edition.

Lang: Calculus of Several Variables. Third edition.

Lang: Introduction to Linear Algebra. Second edition.

Lang: Linear Algebra. Third edition.

Lang: Undergraduate Algebra. Second edition.

Lang: Undergraduate Analysis.

Lax/Burstein/Lax: Calculus with Applications and Computing. Volume 1.

LeCuyer: College Mathematics with APL.

Lidl/Pilz: Applied Abstract Algebra.

Macki-Strauss: Introduction to Optimal Control Theory.

Malitz: Introduction to Mathematical Logic.

(continued)

Undergraduate Texts in Mathematics

Marsden/Weinstein: Calculus I, II, III.
Second edition.

Martin: The Foundations of Geometry
and the Non-Euclidean Plane.

Martin: Transformation Geometry: An
Introduction to Symmetry.

Millman/Parker: Geometry: A Metric
Approach with Models. Second
edition.

Moschovakis: Notes on Set Theory.

Owen: A First Course in the
Mathematical Foundations of
Thermodynamics.

Palka: An Introduction to Complex
Function Theory.

Pedrick: A First Course in Analysis.

Peressini/Sullivan/Uhl: The Mathematics
of Nonlinear Programming.

Prenowitz/Jantosciak: Join Geometries.

Priestley: Calculus: An Historical
Approach.

Protter/Morrey: A First Course in Real
Analysis. Second edition.

Protter/Morrey: Intermediate Calculus.
Second edition.

Ross: Elementary Analysis: The Theory
of Calculus.

Samuel: Projective Geometry.
Readings in Mathematics.

Scharlau/Opolka: From Fermat to
Minkowski.

Sigler: Algebra.

Silverman/Tate: Rational Points on
Elliptic Curves.

Simmonds: A Brief on Tensor Analysis.
Second edition.

Singer/Thorpe: Lecture Notes on
Elementary Topology and Geometry.

Smith: Linear Algebra. Second edition.

Smith: Primer of Modern Analysis.
Second edition.

Stanton/White: Constructive
Combinatorics.

Stillwell: Elements of Algebra:
Geometry, Numbers, Equations.

Stillwell: Mathematics and Its History.

Strayer: Linear Programming and Its
Applications.

Thorpe: Elementary Topics in
Differential Geometry.

Troutman: Variational Calculus and
Optimal Control with Elementary
Convexity. Second edition.

Valenza: Linear Algebra: An
Introduction to Abstract Mathematics.

Whyburn/Duda: Dynamic Topology.

Wilson: Much Ado About Calculus.

This is a textbook on the history, philosophy, and foundations of mathematics. One of its aims is to present some interesting mathematics, not normally taught in other courses, in a historical and philosophical setting. The book is intended mainly for undergraduate mathematics students, but is also suitable for students in the sciences, humanities, and education with a strong interest in mathematics. It proceeds in historical order from about 1800 BC to 1800 AD and then presents some selected topics of foundational interest from the 19th and 20th centuries. Among other material in the first part, the authors discuss the renaissance method for solving cubic and quartic equations and give rigorous elementary proofs that certain geometrical problems posed by the ancient Greeks (e.g. the problem of trisecting an arbitrary angle) cannot be solved by ruler and compass constructions. In the second part, they sketch a proof of Gödel's incompleteness theorem and discuss some of its implications, and also present the elements of category theory, among other topics. The authors' approach to a number of these matters is new.

ISBN 0-387-94544-X

