

This book is a part of the

ALLYN AND BACON SERIES IN ADVANCED MATHEMATICS

**Consulting Editors: Irving Kaplansky
Charles De Prima**

Other books in the series:

James Dugundji, *Topology*

Horst Herrlich and George E. Strecker, *Category Theory*

Irving Kaplansky, *Commutative Rings*

Irving Kaplansky, *Linear Algebra and Geometry: A Second Course*

Ralph Kopperman, *Model Theory and Its Applications*

Joseph J. Rotman, *The Theory of Groups: An Introduction*

Set Theory

AND

Metric Spaces

IRVING KAPLANSKY

University of Chicago

Allyn and Bacon, Inc.

Boston



© Copyright 1972 by Allyn and Bacon, Inc.
470 Atlantic Avenue, Boston

All rights reserved. No part of the material protected by this copyright notice may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying, recording, or by any informational storage and retrieval system, without written permission from the copyright owner.

Library of Congress Catalog Card Number: 71-182353

Printed in the United States of America

to my mother

Contents

PREFACE xi

1

BASIC SET THEORY 1

- 1.1 Inclusion* 1
- 1.2 Operations on Sets* 4
- 1.3 Partially Ordered Sets and Lattices* 9
- 1.4 Functions* 14
- 1.5 Relations; Cartesian Products* 19

2

CARDINAL NUMBERS 21

- 2.1 *Countable Sets* 21
- 2.2 *Cardinal Numbers* 27
- 2.3 *Comparison of Cardinal Numbers; Zorn's Lemma* 31
- 2.4 *Cardinal Addition* 40
- 2.5 *Cardinal Multiplication* 41
- 2.6 *Cardinal Exponentiation* 44

3

WELL-ORDERING; THE AXIOM OF CHOICE 49

- 3.1 *Well-ordered Sets* 49
- 3.2 *Ordinal Numbers* 55
- 3.3 *The Axiom of Choice* 58
- 3.4 *The Continuum Problem* 64

4

BASIC PROPERTIES OF METRIC SPACES 67

- 4.1 *Definitions and Examples* 67
- 4.2 *Open Sets* 71
- 4.3 *Convergence; Closed Sets* 75
- 4.4 *Continuity* 80

5

COMPLETENESS, SEPARABILITY, AND COMPACTNESS 84

5.1 *Completeness* 84

5.2 *Separability* 94

5.3 *Compactness* 98

6

ADDITIONAL TOPICS 106

6.1 *Product Spaces* 106

6.2 *A Fixed-point Theorem* 108

6.3 *Category* 111

APPENDIXES 115

1. *Examples of Metric Spaces* 117

2. *Set Theory and Algebra* 122

3. *The Transition to Topological Spaces* 127

SELECTED BIBLIOGRAPHY 133

INDEX 137

Preface

I first taught a course on set theory and metric spaces in the autumn of 1949. In subsequent years Edwin Spanier presented the material in a somewhat similar way, and he prepared an excellent set of mimeographed notes. These notes were used repeatedly as a text at Chicago.

I have now put them into a somewhat more definitive form. I am very grateful to Spanier for courteously allowing me to complete a project in which he was so deeply involved, and for permission to incorporate numerous exercises from his notes.

The two halves of the book are of nearly equal size. The set theory (with a bow to Halmos) is super-naive. Axiomatic set theory is barely mentioned. The paradoxes get some attention, but in effect they are brushed aside as not really being menacing. My intention is to present the set theory that a working mathematician really needs ninety-nine per cent of the time, and a little bit more (on the theory that to be sure of doing enough, you must do more than enough). A little knowledge may not be a dangerous thing, but a little axiomatic set theory is just not much fun; I hold that one should either take a healthy bite or leave it out in toto.

On the other hand, I hope that this book will help to fight those who say that set theory is a luxury. Hilbert vowed that no one would ever drive us out of the paradise created by Cantor. Let those who agree

strive diligently to transmit to the next generation the knowledge that there is such a paradise.

In the metric space half of the book I have tried to cover the basic topics with a helpful amount of detail and motivation. I hope it will be found useful by teachers who share my belief that topology is best introduced first in the less austere setting of metric spaces. A final appendix has been added to help bridge the gap between metric and topological spaces.

The asterisk attached to occasional exercises probably needs no explanation. I was tempted to put a double asterisk on some. To the teacher: Assign them cautiously!

I am very grateful to Rose Banfield, Fred Flowers, and Tere Shuman for their excellent job in typing drafts of the manuscript. Among the many students who made helpful comments, I wish particularly to thank Susan Bolotin Friedman and Andrew Gallant. To Judith Fiske and Carl Harris of Allyn and Bacon: my gratitude for their contribution to producing the book.

IRVING KAPLANSKY

Set Theory and Metric Spaces

1

Basic Set Theory

1.1 INCLUSION

The word “set” is short and innocuous. In this book we shall treat it as an everyday word, whose meaning everybody knows. For verbal variety we occasionally use a synonym such as “class” or “collection.” This is particularly tempting when we are speaking about a set of sets.

The important thing to know about a set is: What are its members? A handy notation is also a good idea. Given a set A , and a member x of A , we record that information in the notation “ $x \in A$ ”. Officially this is to be read “ x is a member of A ”, but no doubt before long you and I will be reading it as “ x epsilon A ”.

There is also a pleasant notation for describing a set by enumerating its members within a pair of braces. Here are some numerical examples: $\{1, 2, 3, 4, 5\}$, $\{-1, 0, 1\}$, $\{1, 4, 9, 16\}$, $\{\frac{1}{2}, \frac{3}{2}\}$. For instance, the first of these sets has five members, to wit the numbers 1, 2, 3, 4, and 5. We offer two non-numerical examples of sets: $\{\text{Romeo, Juliet}\}$, $\{\text{Hitler, Mussolini, Stalin}\}$.

It may be tedious, difficult, or impossible to list all the members of a set. In that case a suitable use of dots can come to the rescue, as in the examples $\{1, 2, 3, \dots, 98, 99\}$ or $\{1, 2, 3, \dots\}$. In the latter example the

dots replace an infinite number of elements, and the set within the braces is meant to be the set of all positive integers.

A set A may be described by giving a criterion for membership. Example: Let A be the set of all integers x from 1 to 9 satisfying $x^2 + 5x = 14$. The set A can be more closely identified by testing each of the numbers from 1 to 9 in turn. It would be more efficient to solve the equation and then keep only those solutions which are integers from 1 to 9.

We admit the possibility that a set may have no members at all. This memberless set is appropriately called the *empty set*, and the usual notation for it is the symbol \emptyset . A popular alternate name is the *null set*.

Do we really need the somewhat bizarre empty set? The answer is that we could get along without it, but it is very convenient to have the empty set around. For instance, in the next section we shall introduce the intersection of two sets, and it is highly desirable that it be defined for *any* two sets. For this to be the case when the two sets are disjoint (i.e. have no elements in common), we need to be able to say that the intersection is \emptyset .

The number 0 plays a somewhat similar role relative to the positive integers. For centuries, humanity got along reasonably well with just the positive integers, but there was a distinct improvement in the efficiency of arithmetic when the number 0 was invented.

We insert at this point an often repeated warning: Do not confuse the set $\{x\}$, whose only member is x , with x itself. (However, the danger that such confusion will cause serious trouble may be exaggerated.) The possibility of confusion increases if x itself is a set. As a semantic exercise to illustrate the point, consider the sets \emptyset , $\{\emptyset\}$, and $\{\{\emptyset\}\}$. The first has no members; the second has exactly one member, namely \emptyset ; and the third has as its sole member the set whose only member is \emptyset .

To conclude this section we discuss the inclusion of one set in another, and the equality of two sets.

We say that A is *included* in B , and write $A \subset B$, if every member of A is a member of B . Equally popular is the terminology " A is contained in B ". Equality of A and B is allowed; some authors emphasize this by using the notation " \subseteq " or " \supseteq " instead. It is on the theory that the most important concept should get the simplest notation that we are using " \subset ". On the infrequent occasions when we must carefully note that equality is excluded, we shall either put it verbally (A is properly included in B), or use the (admittedly clumsy) notation $A \subsetneq B$.

In the preceding paragraph we took equality of sets as self-explanatory; of course two sets are equal precisely when they have the same members. Worth adding is the following equally obvious fact: If A and B are sets satisfying $A \subset B$ and $B \subset A$, then $A = B$. Why bother to mention this at all? One answer is that it is helpful in proving that two sets are equal. Of

course, if we are able to see at a glance that $x \in A$ holds if and only if $x \in B$ does, then $A = B$ has been forthwith established. But in a situation of sufficient complexity it may be wise strategy to separate the problem into the two halves, $A \subset B$ and $B \subset A$. Such a proof looks as follows: We first show that $x \in A$ implies $x \in B$ and then that $x \in B$ implies $x \in A$. Later we shall see several examples of proofs in this style.

We record the "transitivity" of inclusion: If $A \subset B$ and $B \subset C$, then $A \subset C$. We dismiss this also as too trivial for a formal proof.

If $B \subset A$ we say that B is a *subset* of A . Observe that any set A has the obvious subsets \emptyset and A . The set of all subsets of A is called the *power set* of A . Our notation for it is $P(A)$.

On occasion we shall reverse the inclusion symbol: $A \supset B$ (read " A includes B " or " A contains B ") is equivalent to $B \subset A$. Not entirely consistently, we also say that A contains x when x is a member of A instead of a subset of A .

EXERCISES

1. A set A has the property that $A \subset B$ holds for any set B . Prove that $A = \emptyset$.
2. Let A , B , and C be sets satisfying $A \subset B$, $B \subset C$, and $C \subset A$. Prove that $A = B = C$.
3. How many elements are there in the following sets: \emptyset , $\{\emptyset\}$, $\{\{\emptyset\}\}$, $\{\emptyset, \{\emptyset\}\}$, $\{\emptyset, \emptyset\}$?
4. List all inclusions that hold among the following sets:
 - (a) $A = \{2, 4, 6\}$,
 - (b) $B = \{2, 4, 6, 8\}$,
 - (c) $C = \emptyset$,
 - (d) $D =$ all even integers between 1 and 9.
5. List all inclusions that hold among the following sets:
 - (a) $A =$ all integers from 1 to 9 satisfying $x^2 - 5x = 14$,
 - (b) $B = \{2, 7\}$,
 - (c) $C = \{-2, 7\}$,
 - (d) $D = \{7\}$.
6. List all the subsets of $\{1, 2\}$. How many subsets are there?
7. List all the subsets of $\{1, 2, 3\}$. How many subsets are there?
8. The preceding two exercises have invited the guess that a set with n elements has 2^n subsets. Prove this. (*Hint:* You can argue directly that n decisions have to be made, in each of which there are two possibilities: including or excluding a given element. Alternatively, the proof may be given by induction. Assume that there are 2^{n-1} subsets of $\{1, 2, \dots, n-1\}$. To each of these n might or might not be adjoined.)

1.2 OPERATIONS ON SETS

The fundamental operations on sets are intersection and union.

The *intersection* of two sets A and B is the set of all elements lying in both A and B . Our notation for it is $A \cap B$ (suggested reading: “ A meet B ”). The characteristic property of $A \cap B$ can be stated as follows: $x \in A \cap B$ if and only if $x \in A$ and $x \in B$.

We offer three examples.

- (a) If $A = \{2, 4\}$ and $B = \{2, 7\}$, then $A \cap B = \{2\}$.
- (b) If $A = \{2, 3, 7\}$ and $B = \{2, 7\}$, then $A \cap B = \{2, 7\}$.
- (c) If $A = \{2, 3\}$ and $B = \{4, 7\}$, then $A \cap B = \emptyset$.

We proceed to examine the operation of intersection to see whether it obeys some of the “laws of algebra.” The commutative and associative laws in fact hold:

- (1) $A \cap B = B \cap A.$
- (2) $(A \cap B) \cap C = A \cap (B \cap C).$

Equations (1) and (2) are too obvious to merit formal proof. So are the next two statements, which link intersection with inclusion and with the null set:

- (3) $A \cap B = A$ if and only if $A \subset B.$
- (4) For any A , $A \cap \emptyset = \emptyset.$

The *union* of two sets A and B is the set of all elements lying in A or B . (The vagaries of the English language offer a possible pitfall, so we state more exactly that we mean elements lying in A or B or *both*. It is actually interesting—see Exercise 5—to examine the possibility of taking the elements lying in A or B but not in both.) We write $A \cup B$ (read “ A union B ”) for the union of A and B .

In the three examples above, $A \cup B$ works out to be $\{2, 4, 7\}$ in (a), $\{2, 3, 7\}$ in (b), and $\{2, 3, 4, 7\}$ in (c).

Analogues of (1)–(4) are valid.

- (1') $A \cup B = B \cup A.$
- (2') $(A \cup B) \cup C = A \cup (B \cup C).$
- (3') $A \cup B = A$ if and only if $B \subset A.$
- (4') For any A , $A \cup \emptyset = A.$

There is a pictorial representation of union and intersection that can be helpful. We show A and B as more or less circular objects, and place them in more or less general positions. In Figure 1, $A \cap B$ is shaded; in

Figure 2, $A \cup B$ is shaded. Such figures are called *Venn diagrams*. We must be scrupulously careful to use Venn diagrams only for motivation, not as a substitute for a proof.

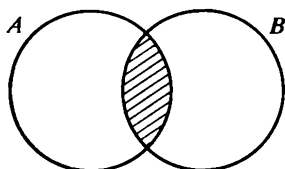
 $A \cap B$

Figure 1

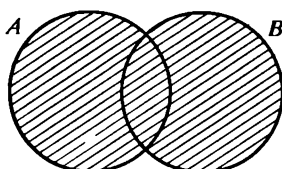
 $A \cup B$

Figure 2

John Venn (1834–1923), an English logician, discussed such pictorial representations in *Proc. Camb. Phil. Soc.* 4 (1880), 36–59.

It is tempting to pursue an analogy between the set-theoretic operations of intersection and union and the operations of addition and multiplication on ordinary numbers. If, for instance, we experiment with \cap , \cup , and \emptyset by treating them as multiplication, addition, and 0, respectively, we find that (1), (2), (4) and (1'), (2'), (4') stand up splendidly. This might well make us bold enough to try

$$(5) \quad A \cap (B \cup C) = (A \cap B) \cup (A \cap C),$$

the analogue of the distributive law $a(b + c) = ab + ac$. Venn diagrams representing $A \cap B$, $A \cap C$, and $B \cup C$ (Figure 3) are encouraging. The reader should proceed to visualize the two sides of (5) in Figure 3, to see if they look equal.

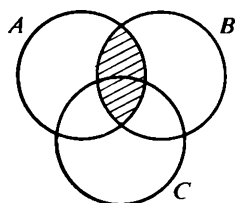
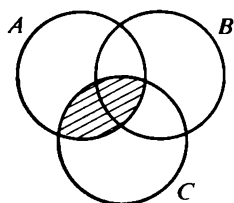
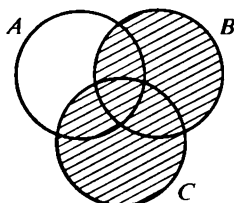
 $A \cap B$  $A \cap C$  $B \cup C$

Figure 3

We give a formal proof of (5), in the style mentioned at the end of Section 1.1.

(a) Suppose that x lies in $A \cap (B \cup C)$, the left side of (5). Then x lies in both A and $B \cup C$. The statement $x \in B \cup C$ implies that x lies in B or C or both. Assume that $x \in B$. Then x lies in both A and B , $x \in A \cap B$, and hence x lies in $(A \cap B) \cup (A \cap C)$, the right side of (5). If $x \in C$, then $x \in A \cap C$, and again $x \in (A \cap B) \cup (A \cap C)$.

(The cases $x \in B$ and $x \in C$ actually need not have been distinguished, for there is perfect symmetry between B and C . In the second half of the proof we shall take advantage of this symmetry.)

(b) Suppose that x lies in $(A \cap B) \cup (A \cap C)$, the right side of (5). Then $x \in A \cap B$ or $x \in A \cap C$. By symmetry we may assume that $x \in A \cap B$. This gives us both $x \in A$ and $x \in B$. The statement $x \in B$ implies the weaker statement $x \in B \cup C$. Hence x lies in both A and $B \cup C$, i.e. $x \in A \cap (B \cup C)$, the left side of (5).

Since the statements (1)–(4) were accompanied by (1')–(4'), in which the roles of \cap and \cup got interchanged, it is a reasonable gambit to try

$$(5') \quad A \cup (B \cap C) = (A \cup B) \cap (A \cup C).$$

But the analogous statement for addition and multiplication is the dubious looking equation $a + bc = (a + b)(a + c)$. Although the analogy has broken down, equation (5') is true; we leave it as Exercise 1.

It is a very interesting fact that the analogy between the set-theoretic operations and addition and multiplication can be fully restored by replacing union by the *symmetric difference* defined in Exercise 9. This discovery is due to M. H. Stone.

The pioneering work of Marshall H. Stone (1903–) on Boolean algebras included the discovery (just mentioned) that they could be transformed into Boolean rings, in the sense of ordinary ring theory. Other important achievements included the Stone-Cech compactification and the Weierstrass-Stone approximation theorem. He is probably best known for his authoritative treatise on Hilbert space. The major part of his career was spent at Harvard University and at the University of Chicago.

We introduce an additional operation: complementation. In order to do so, we have to make an agreement that everything is happening inside some big fixed set. Call such a set U (the letter is meant to suggest "universal"). Given a set A with $A \subset U$ we define A' , the *complement* of A , to be the set of all elements lying in U but not in A . Then we have

$$(6) \quad (A')' = A,$$

$$(7) \quad A \cap A' = \emptyset, \quad A \cup A' = U.$$

Furthermore, for A and B both included in U :

$$(8) \quad (A \cap B)' = A' \cup B',$$

$$(9) \quad (A \cup B)' = A' \cap B'.$$

The proof of (6)–(9) is left as Exercise 2.

Complementation provides us with a systematic way of interchanging the roles of union and intersection. (This is described as a *principle of duality* for \cup and \cap .) Rather than give complete details, we shall illustrate the procedure by showing how to pass from (5) to (5'). (Nevertheless we urge the reader to give a direct proof of (5') in Exercise 1.)

We make a preliminary remark. In order to prove two sets (say X and Y) equal it is sufficient to prove that their complements X' and Y' are equal. For we apply complementation to the equation $X' = Y'$, getting $(X')' = (Y')'$; by (6) this yields $X = Y$.

We turn to the problem of deriving (5') from (5). For the universal set U (within which we take complements) we can choose $U = A \cup B \cup C$. By the remark in the preceding paragraph it will suffice to prove

$$(5'') \quad [A \cup (B \cap C)]' = [(A \cup B) \cap (A \cup C)]'.$$

Write LS , RS for the left and right sides of (5''). Using (9) and then (8) we obtain

$$LS = A' \cap (B \cap C)' = A' \cap (B' \cup C').$$

We can now use equation (5), with the sets A , B , C replaced by A' , B' , C' . The result is

$$LS = (A' \cap B') \cup (A' \cap C').$$

We next apply (8) and (9) in succession to RS :

$$RS = (A \cup B)' \cup (A \cup C)' = (A' \cap B') \cup (A' \cap C').$$

We have identified LS and RS .

A more general concept than complementation is that of *subtraction* of sets. Given any two sets A and B we define $A - B$ to be the set of all x such that $x \in A$ and $x \notin B$ (in words: elements lying in A but not in B). It is not required that B be a subset of A ; however, we can arrange this at our pleasure by replacing $A - B$ by $A - (A \cap B)$. Note that the complement A' , in the discussion above, is $U - A$ in the notation of subtraction. We have left the main properties of subtraction as Exercise 7.

To conclude this section we discuss the intersection and union of more than two sets. Given three sets A , B , and C we wish to define their intersection. We could be very formal. Having treated the intersection of two sets, we could observe that the two sides of (2) are equal and declare their common value to be the intersection of A , B , and C . The procedure could then be extended inductively to n sets. Besides being somewhat heavy-handed for the intersection of a finite number of sets, this procedure is problematic for an infinite number of sets. We therefore make a fresh start and straightforwardly define the intersection $A \cap B \cap C$ of A , B , and C to be the set of elements lying in all of A , B , and C . At our pleasure we can note that the intersection can be regrouped as either of the terms in (2).

For sets A_1, A_2, \dots, A_n no novel thought is needed; $A_1 \cap A_2 \cap \dots \cap A_n$

A_n is the set of all x such that $x \in A_i$ for all i . Nor is there any real difficulty with the intersection of an infinite number of sets; our only problem is to agree on notation. It is inadvisable to write something like $A_1 \cap A_2 \cap A_3 \dots$, for this would suggest that the number of sets is countable (we are anticipating, for the moment, the discussion of countability in Section 2.1). What we need is the concept of a general *index set*. Let I be any arbitrary set; we write i for a typical member of I . Let a set A_i be given for each $i \in I$. The intersection of all of the sets A_i , written $\bigcap_i A_i$, is the set of all x satisfying $x \in A_i$ for all $i \in I$. The union $\bigcup_i A_i$ is treated analogously.

Some properties of these unlimited unions and intersections are assembled in Exercise 8.

EXERCISES

1. Prove (5').
2. Prove (6), (7), (8), and (9).
3. Let N be the set of all positive integers, A the set of even integers, B the set of odd integers, and C the set of multiples of three.
 - (a) Describe the sets $A \cap C$, $B \cap C$, $B \cup C$. (C' is the complement of C within N .)
 - (b) Verify that $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$.
4. Prove that $(A \cap B) \cup C = A \cap (B \cup C)$ if and only if $C \subset A$.
5. If A and B are subsets of U , and B' denotes the complement of B within U , prove that $A \cap B' = \emptyset$ if and only if $A \subset B$.
6. For a finite set A , let $\alpha(A)$ denote the number of elements in A . If A and B are finite sets, prove that

$$\alpha(A \cap B) + \alpha(A \cup B) = \alpha(A) + \alpha(B).$$

7. For any sets A , B , and C , prove:
 - (a) $C - (A \cup B) = (C - A) \cap (C - B)$.
 - (b) $C - (A \cap B) = (C - A) \cup (C - B)$.
 - (c) $B - (B - A) = A \cap B$.
 - (d) $(A - B) \cup (B - A) = (A \cup B) - (A \cap B)$.

(An alternative to proving these from scratch is to reduce them to properties of intersection, union, and complementation by observing that $A - B = A \cap B'$.)

8. With an arbitrary (possibly infinite) index set I , prove:
 - (a) $A \cup (\bigcap_i B_i) = \bigcap_i (A \cup B_i)$.
 - (b) $A \cap (\bigcup_i B_i) = \bigcup_i (A \cap B_i)$.

In (c) and (d) the complementation is taken relative to a fixed universal set U .

$$(c) \quad \left(\bigcup_i A_i\right)' = \bigcap_i A_i'.$$

$$(d) \quad \left(\bigcap_i A_i\right)' = \bigcup_i A_i'.$$

9. The *symmetric difference* of sets A and B is $(A - B) \cup (B - A)$, i.e., all elements in A or in B but not in both. Write $A + B$ for the symmetric difference and shorten $A \cap B$ to AB .
- (a) Prove that $A(B + C) = AB + AC$.
- (b) If you know the relevant definitions, prove that these operations make $P(A)$ a commutative associative ring with unit in which every element is idempotent. (Recall that $P(A)$ is the power set of A —the set of all subsets of A .)
10. Call a subset B of a set A *cofinite* if the complement of B in A is finite. If B and C are cofinite subsets of A , prove that $B \cap C$ is cofinite.

1.3 PARTIALLY ORDERED SETS AND LATTICES

This section has two purposes. We wish to insert set-theoretic inclusion and the operations on sets into a broader context, in order to get a better perspective. The broader context will also play a key role in our later discussion of Zorn's lemma and the axiom of choice.

DEFINITION. Let L be a set. We work with a symbol \leq ("less than or equal to"). Given elements a and b in L , the statement $a \leq b$ may or may not be valid, depending on a and b . We impose three axioms.

1. For all a in L , $a \leq a$.
2. If $a \leq b$ and $b \leq a$, then $a = b$.
3. If $a \leq b$ and $b \leq c$, then $a \leq c$.

When these axioms are satisfied, we say that L is a *partially ordered set*.

Remarks: 1. Instead of $a \leq b$ we may write equivalently $b \geq a$. If $a \leq b$ and $a \neq b$ we write $a < b$ or $b > a$.

2. In the case of set-theoretic inclusion we explained the wisdom of allowing \subset to admit the case of equality. Why not do the same here? The answer: tradition. It is ingrained in us that, for numbers, $a < b$ means that a is strictly smaller than b , and we are now imitating numbers rather than sets.

We define at once an important special class of partially ordered sets.

DEFINITION. A *chain* is a partially ordered set in which, for any a and b , either $a \leq b$ or $b \leq a$. (Other names: linearly ordered set, simply ordered set.)

We present a number of examples of partially ordered sets.

(1) Assume that, in L , $a \leq b$ holds only when $a = b$. Then L is a partially ordered set. It is suggestive to call L "totally unordered." In a sense the totally unordered case is at the opposite extreme from the case of a chain.

(2) The set of all real numbers forms a chain in its natural ordering.

(3) Let L be a partially ordered set and M a subset of L . If we say that $a \leq b$ is to hold for $a, b \in M$ precisely when it does in L , we make M into a partially ordered set. If L is a chain, so is M . Thus any set of real numbers is a chain in its natural ordering. Possible subsets of the chain of real numbers: the rational numbers, the integers, the positive integers.

(4) Let L be a finite chain. Evidently L has a smallest element, then a next smallest, etc. We conclude that if a chain has a finite number of elements, say n , it looks exactly like the set $\{1, 2, \dots, n\}$ in its natural order $1 < 2 < \dots < n$.

(5) Partially ordered sets with a small number of elements can be surveyed by inspection. If L is a partially ordered set with 2 elements, then L is either totally unordered or a chain. If L has 3 elements, there are 5 possibilities. They are most easily described by the diagrams exhibited in Figure 4, where we draw a line going up from a to b if $a < b$ and b is directly above a (i.e. there is no c with $a < c < b$).

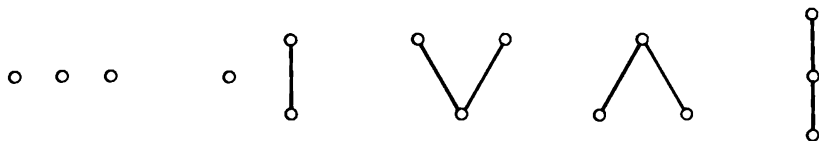


Figure 4

(6) Let A be a fixed set and $P(A)$ its power set (the set of all subsets of A). By using set-theoretic inclusion \subset for \leq we make $P(A)$ into a partially ordered set. Figure 5 illustrates $P(A)$ for $A = \{1, 2\}$ and $A = \{1, 2, 3\}$.

We seek to distinguish the partially ordered sets $P(A)$ of Example 6 among all partially ordered sets. The notions of upper bound and least upper bound play an important role here, and will also be used in many later contexts.

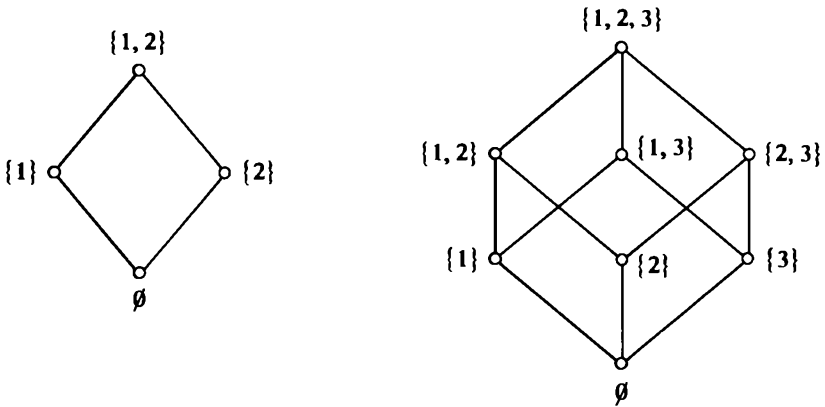


Figure 5

DEFINITIONS. Let S be a subset of a partially ordered set L . An element $u \in L$ is called an *upper bound* of S if $s \leq u$ for all $s \in S$. We say that u is a *least upper bound* of S if it is an upper bound of S and furthermore $u \leq v$ for any upper bound v of S . *Lower bound* and *greatest lower bound* are defined analogously.

It is immediate that if a least upper bound for a subset exists, then it is unique. Thus it is appropriate to speak of “*the* least upper bound,” and analogously “*the* greatest lower bound.”

It is important to note that an upper bound or least upper bound for a subset S need not lie in S . To emphasize the distinction, we introduce the terms “*top element*” and “*bottom element*.” An element u of a partially ordered set L is a *top element* of L if $a \leq u$ holds for all a in L . *Bottom element* is defined analogously. It is a severe requirement on a partially ordered set to demand that every subset have a top element and a bottom element—see Exercise 2.

A subset S of a partially ordered set L need not have an upper bound. For instance, we can take S to be the entire set L , and an upper bound for all of L exists if and only if L has a top element. Thus if L is the chain of all real numbers, L does not have an upper bound. This can be remedied by supplying a new element (call it ∞) sitting on top of the chain of real numbers. If we similarly augment the real numbers at the bottom by inserting $-\infty$, we get a chain, called the extended real numbers, which has the property that every subset has a least upper bound and greatest lower bound.

As a second example, take L to be the chain of rational numbers x satisfying $0 \leq x \leq 2$. Here there is a top element and a bottom element,

so every subset has both an upper bound and a lower bound. However, a least upper bound need not exist: take S to be the subset of L consisting of those x with $x^2 < 2$.

A finite example (necessarily not a chain) in which a least upper bound fails to exist is provided by the partially ordered set shown in Figure 6. Here c and d are upper bounds for $S = \{a, b\}$, but S has no least upper bound.

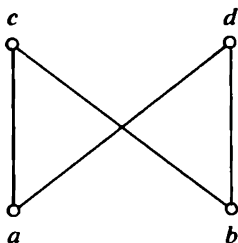


Figure 6

DEFINITION. A *lattice* is a partially ordered set in which every two elements have a least upper bound and a greatest lower bound. The notation is $a \cup b$ for the least upper bound and $a \cap b$ for the greatest lower bound.

We readily see that in a lattice any finite subset has a least upper bound and a greatest lower bound. This need not be true for infinite subsets. We say that a lattice L is *complete* if any subset of L has a least upper bound and a greatest lower bound.

The equations (5) and (5') in Section 1.2 make sense in any lattice, but they need not hold (see Exercise 8). We say that a lattice is *distributive* if they do hold. (It can be proved that we need not assume both (5) and (5'); either one implies the other.)

Let L be a lattice possessing a top element and a bottom element, for which we use the notation 1 and 0 , respectively. Thus we have $0 \leq x \leq 1$ for every x in L . We say that L is *complemented* if for any a in L there exists b in L with $a \cap b = 0$, $a \cup b = 1$. (We might have written a' instead of b , but we did not wish to suggest that the complement is to be unique, in the way that it is when L is the power set of a set. In this connection see Exercise 9.)

We can summarize Section 1.2 as follows: For any set A , consider its power set $P(A)$ as a partially ordered set under inclusion. Then $P(A)$ is a complete complemented distributive lattice.

There is a partial converse. We put aside the question of completeness, and collect the remaining properties in a definition.

DEFINITION. A *Boolean algebra* is a complemented distributive lattice.

George Boole (1815–1864), in his treatise *An Investigation of the Laws of Thought* (1854), devised a systematic symbolic way of dealing with logic. It is equally appropriate to the algebra of sets, as developed in Section 1.2.

If L is a finite Boolean algebra, then it can be proved that there exists a finite set A such that L is a perfect replica of $P(A)$. (In the terminology to be introduced in the next section, L is isomorphic to $P(A)$.) This is not true in the infinite case, but the weaker statement can be made that, for a suitable A , L is identifiable with a subset of $P(A)$; more exactly, L is identifiable with a subcollection of the subsets of A , where the subcollection contains \emptyset and A and is closed under intersections, unions, and complements. The proofs of these statements are not difficult, but a discussion would take us too far out of our way (and the infinite case requires the transfinite methods we have still to develop).

EXERCISES

1. Let a symbol $<$ be given for a set L . Assume that $a < a$ is never true and that $a < b$ and $b < c$ imply $a < c$. Define $a \leq b$ to mean that $a = b$ or $a < b$. Prove that L is a partially ordered set relative to \leq .
- 2.* Let L be a partially ordered set in which every subset has a top and bottom element. Prove that L is a finite chain.
3. Let N be the chain of positive integers, in its usual order. Is N complete? Is N complete if infinity is placed on top?
4. Let N be the set of positive integers and define $m \leq n$ to mean that m divides n . Is N a lattice? Is it complete?
5. Let A be a set such that its power set $P(A)$ is a chain (with inclusion as the partial ordering). What can be said about A ?
6. Prove that any chain is a distributive lattice.
7. Draw figures for all partially ordered sets with 4 elements. How many are lattices?
8. Draw figures for the 5 different lattices with 5 elements. Find 2 of these which are not distributive.
9. Let L be a distributive lattice with a top element and a bottom element. Prove: If an element of L has a complement, the complement is unique.

10. Let L be a partially ordered set in which every subset has a least upper bound. Suppose that L has a bottom element. Prove that L is a complete lattice (i.e. every subset has a greatest lower bound). (*Hint*: Given a subset S , look at the set of lower bounds of S .)
11. Let L be a lattice in which every subset with an upper bound has a least upper bound. Prove that any subset of L with a lower bound has a greatest lower bound. (A lattice with either—hence both—of these properties is called *conditionally complete*.)
- 12.* Call a subset S of a chain L a *lower segment* if it has the following property: If $x \in S$ and $a < x$, then $a \in S$. Prove the equivalence of the following two statements: (a) L is conditionally complete, (b) For any lower segment S of L other than L and \emptyset , there exists an element u such that S is either the set of all x with $x \leq u$, or the set of all x with $x < u$.
- 13.* Call a subset S of a chain *convex* if it has the following property: If $a, b \in S$ and $a < x < b$, then $x \in S$. Prove that a convex subset of a chain L is the intersection of a lower segment of L and an upper segment of L (upper segments being defined in the same way as the lower segments of the preceding exercise).
- 14.* Prove that an infinite partially ordered set contains either an infinite chain or an infinite totally unordered subset.
- 15.* Let L be a partially ordered set in which the maximum length of a subchain is n (n finite). Prove that L is the union of n totally unordered subsets and that n is the smallest integer with this property.
- 16.* Let L be a finite partially ordered set in which the maximum size of a totally unordered subset is n . Prove that L can be expressed as a union of n subchains. (*Remark*: Although this is quite similar to the preceding exercise, it is harder. Finiteness is not necessary, but removing it requires transfinite methods.)

1.4 FUNCTIONS

We use the term “function” in the ordinary sense familiar from elementary mathematics. A *function* f from A to B assigns to each a in A an element $f(a)$ in B . We describe the setup by the symbol $f: A \rightarrow B$. We call A the *domain* of f ; note that f is defined on all of A . The *range* of f is the set of all elements $f(a)$, with a any element in A ; the range may or may not be all of B . We sometimes refer to the range as the *image* of A under f , and an individual element $f(a)$ as the image of a .

Alternative terms for “function” are “map” or “mapping.”

We say that f is *one-to-one* if $f(a_1) = f(a_2)$ implies $a_1 = a_2$ for all a_1, a_2 in A . We say that f is *onto* if every b in B is the image of some element of A . In wide use are the adjectives “injective” instead of “one-to-one” and

“surjective” instead of “onto.” In case of any doubt, it is wise to say “onto B ”; observe that every function maps onto its range.

Crude pictures are helpful in visualizing functions. We exhibit A and B as irregular cloud-shaped objects and exhibit arrows running from A to B ; they are meant to suggest the act of picking up points of A and moving them to B . In Figure 7 the function illustrated fails to be one-to-one since two points have the same image; also the range fails to fill up B .

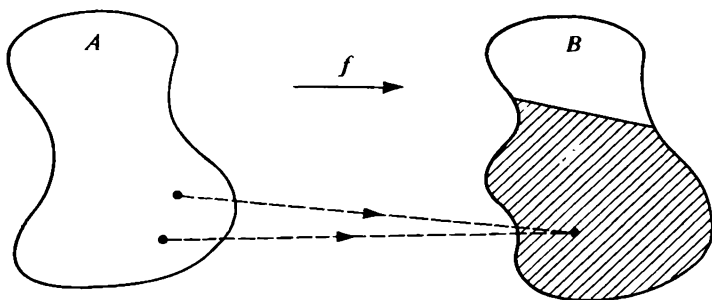


Figure 7

A function which is both one-to-one and onto is a *one-to-one correspondence*. Of special importance is the one-to-one correspondence $i: A \rightarrow A$ defined by $i(a) = a$ for all a . We call it the *identity function*. More exactly, it is the identity function on A , and in case of ambiguity the symbol i_A may be used.

Let $f: A \rightarrow B$ and $g: B \rightarrow C$ be functions. Then a *product* or *composite* function gf is defined by setting $(gf)(a) = g(f(a))$. The reader should note carefully that to get the product function gf you first apply f , then g .

The associative law holds. To state it in maximum generality we need four sets and three functions, say $f: A \rightarrow B$, $g: B \rightarrow C$, and $h: C \rightarrow D$ (Figure 8).

We dismiss the equation $(hg)f = h(gf)$ as obvious.

It might happen that we have functions f and g from A to B and that

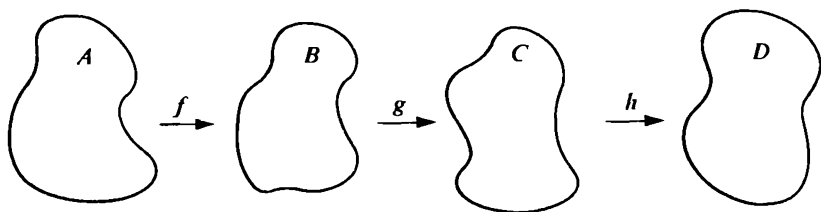


Figure 8

a product is defined on B . (For instance let B be the integers with their usual multiplication.) Then there is another meaning which gf might have: $(gf)(a) = g(a)f(a)$. In case of possible ambiguity it might be wise to write $g(f)$, or perhaps $g \circ f$, for the composite being discussed in this section.

Let f be a one-to-one correspondence from A to all of B . There is a natural definition of a function f^{-1} from B to A : we define $f^{-1}(b)$ to be the unique a such that $f(a) = b$. In the pictorial representation we run the arrow back from B to A . Observe that $f^{-1}f$ is the identity function on A and ff^{-1} is the identity function on B .

If $f: A \rightarrow B$ is not one-to-one, an attempt to define f^{-1} is hampered by our inability to pick an inverse image uniquely. If f is one-to-one but not onto, we are unable to define f^{-1} on all of B (we can of course define it on the range of f).

The picture brightens considerably if we give up the attempt to define f^{-1} on elements and define it instead on subsets. Once we decide to do this, we realize that there is an even more obvious function on subsets going in the direction $A \rightarrow B$. Thus, let f be any function from a set A to a set B . For a subset S of A , the obvious meaning of $f(S)$ is the set of all $f(s)$ with $s \in S$. In this way we have defined a map from the power set $P(A)$ of A to the power set $P(B)$. Observe that the symbol we are using for this map is the same letter f —logically objectionable but ever so convenient. For a subset T of B , we define $f^{-1}(T)$ to be the set of all a in A with $f(a) \in T$, and we call $f^{-1}(T)$ the *inverse image* of T . In sum, $f: A \rightarrow B$ induces $f: P(A) \rightarrow P(B)$ and $f^{-1}: P(B) \rightarrow P(A)$. Of special interest is the inverse image of an element $b \in B$. It is a subset of A (the null set if b is not in the range of f). We write it $f^{-1}(b)$, although a purist might insist on $f^{-1}(\{b\})$.

Since $P(A)$ and $P(B)$ are lattices, we can ask how f and f^{-1} behave relative to the lattice operations (intersection, union, complement, subtraction). It turns out that some things are true and some are not, perhaps a little unpredictably. One general observation is that f^{-1} behaves better than f . See Exercises 9–12.

When two sets A and B have additional structure, it is reasonable to impose restrictions on functions from A to B . Consider, for example, the case where A and B are partially ordered sets. By an *isomorphism* of A onto B we then mean a function f which is one-to-one and onto and which in addition satisfies the following: for all a_1, a_2 in A , $f(a_1) \leq f(a_2)$ holds if and only if $a_1 \leq a_2$. Two isomorphic partially ordered sets are essentially identical; only the “names” of the sets and their members have changed.

Remark: For the sake of emphasis and unambiguity, we often call an isomorphism between partially ordered sets an *order-isomorphism* (as in Exercise 19).

EXERCISES

- Let $f: A \rightarrow B$ and $g: B \rightarrow C$ be given functions.
 - If f and g are one-to-one, prove that gf is one-to-one.
 - If gf is one-to-one, prove that f is one-to-one.
 - If f is onto and gf is one-to-one, prove that g is one-to-one.
 - Give an example where gf is one-to-one, but g is not.
- Let $f: A \rightarrow B$ and $g: B \rightarrow C$ be given functions.
 - If f and g are onto, prove that gf is onto.
 - If gf is onto, prove that g is onto.
 - If gf is onto and g is one-to-one, prove that f is onto.
 - Give an example where gf is onto but f is not. (See Exercise 13 for the definitive version of Exercises 1 and 2.)
- Prove that a function $f: A \rightarrow B$ is one-to-one and onto if and only if there exists $g: B \rightarrow A$ with $gf = i_A$ and $fg = i_B$.
- Suppose that $f: A \rightarrow B$ and $g: B \rightarrow C$ are both one-to-one and onto. Prove that gf is one-to-one and onto [compare Exercises 1(a) and 2(a)]. Prove further that $(gf)^{-1} = f^{-1}g^{-1}$.
- Let N be the set of positive integers. Let $f: N \rightarrow N$ and $g: N \rightarrow N$ be defined by $f(x) = x^2$ and $g(x) = x + 1$. Give formulas for gf and fg . Are these two functions equal? How do they compare with the function given by $f(x)g(x)$?
- Let N denote the positive integers, Z all the integers, and R the real numbers. Which of the following functions is one-to-one? Which is onto?
 - $f(x) = x^2$ on N , on Z , on R .
 - $f(x) = x^3$ on N , on Z , on R .
- Given a function $f: A \rightarrow A$, we write f^n for the function on A obtained by taking the composite of f with itself n times. Suppose that f^n equals the identity function for some n (one then says that f is *periodic*). Prove that f is one-to-one and onto.
- As a generalization of periodic functions (see the preceding exercise) we say that $f: A \rightarrow A$ is *locally periodic* if for every $x \in A$ there exists an integer $n(x)$, depending on x , such that $f^{n(x)}(x) = x$. Prove that a locally periodic function is one-to-one and onto.
- Let $f: A \rightarrow B$ be a given function.
 - If S and T are subsets of A with $S \subset T$, prove that $f(S) \subset f(T)$.
 - If S and T are subsets of B with $S \subset T$, prove that $f^{-1}(S) \subset f^{-1}(T)$.
 - Prove that ff^{-1} is the identity on $P(B)$ if f is onto; and that $f^{-1}f$ is the identity on $P(A)$ if f is one-to-one.
- Given a function $f: A \rightarrow B$, subsets S_1, S_2 of A and subsets T_1, T_2 of B , prove the following:
 - $f(S_1 \cup S_2) = f(S_1) \cup f(S_2)$.
 - $f^{-1}(T_1 \cap T_2) = f^{-1}(T_1) \cap f^{-1}(T_2)$.

- (c) $f(S_1 \cap S_2) \subset f(S_1) \cap f(S_2)$.
 (d) $f^{-1}(T_1 \cap T_2) = f^{-1}(T_1) \cap f^{-1}(T_2)$.

In part (c), show that equality holds for all S_1 and S_2 if and only if f is one-to-one.

11. Generalize Exercise 10 to the intersection and union of any number of subsets.
12. Let $f: A \rightarrow B$ be a given function. Complements are taken within A for subsets of A , and within B for subsets of B .
- (a) Prove that $f^{-1}(T^c) = [f^{-1}(T)]^c$ for any $T \subset B$.
 (b) Prove that $f(S^c) \supset f(S)^c$ for any $S \subset A$ if and only if f is onto.
 (c) Prove that $f(S^c) \subset f(S)^c$ for any $S \subset A$ if and only if f is one-to-one.
13. Let $f: A \rightarrow B$ and $g: B \rightarrow C$ be given functions.
- (a) Prove that gf is one-to-one if and only if f is one-to-one and g is one-to-one on the range of f .
 (b) Prove that gf is onto if and only if g is onto on the range of f .
14. We work within a fixed set A . For a subset S of A , the *characteristic function* ϕ_S of S is a function from A to $\{0, 1\}$ which takes the value 1 on S and 0 on S^c . Prove:
- (a) $\phi_{S \cap T} = \phi_S \phi_T$.
 (b) $\phi_{S^c} = 1 - \phi_S$.
 (c) $\phi_S + \phi_T = \phi_{S \cup T} + \phi_{S \cap T}$.
- 15.* Let L be a complete lattice, and let $f: L \rightarrow L$ be a function for which $a \leq b$ implies $f(a) \leq f(b)$. Prove that f leaves some element of L fixed. (*Hint*: Consider all a in L with $f(a) \geq a$.)
16. Let A be a chain and B a partially ordered set. Let $f: A \rightarrow B$ be a one-to-one function for which $a \leq b$ implies $f(a) \leq f(b)$. Prove that $f(a) \leq f(b)$ implies $a \leq b$.
- 17.* Let f be a map of a set A into itself.
- (a) If B and C are subsets with $f(B) = B$ and $f(C) = C$, prove that $f(B \cup C) = B \cup C$.
 (b) Extend part (a) to any number of subsets.
 (c) Prove that there is a unique largest subset of A which is carried onto itself by f .
 (d) If B and C are subsets with f one-to-one on B and C , is f necessarily one-to-one on $B \cup C$?
18. (a) Let x_0 be a fixed element in a set X . Let $f: X \rightarrow X$ be defined by $f(x) = x_0$ for all x . Prove that $fg = f$ for all $g: X \rightarrow X$.
 (b) Conversely assume that $f: X \rightarrow X$ satisfies $fg = f$ for all $g: X \rightarrow X$. Prove that, for a suitable x_0 , f is the map of part (a).
- 19.* Prove that an infinite chain contains either a chain order-isomorphic to the positive integers or a chain order-isomorphic to the negative integers.

1.5 RELATIONS; CARTESIAN PRODUCTS

Since we have placed functions ahead of relations in our exposition, it is tempting to say that a relation is a function which is allowed to be multiple-valued and is not necessarily everywhere defined. A formal definition is based on the concept of an *ordered pair* (a, b) . The crucial property of ordered pairs is the following: $(a, b) = (c, d)$ if and only if $a = c$ and $b = d$. A *relation* between sets A and B is a set of ordered pairs (a, b) with $a \in A$ and $b \in B$.

Many of the concepts discussed for functions in Section 1.4 (e.g. products, inverses) can be defined in the broader context of relations.

Relations will receive only passing attention in this book. Perhaps the most important special cases are: functions, partial orderings, and equivalence relations. We now sketch the theory of equivalence relations.

Let A be a set and let a relation \sim be given on A . This means that for certain ordered pairs (a, b) of elements of A we have $a \sim b$. We say that \sim is an equivalence relation if the following three properties hold:

1. For all a , $a \sim a$.
2. $a \sim b$ implies $b \sim a$.
3. $a \sim b$ and $b \sim c$ imply $a \sim c$.

These properties are called, respectively, *reflexivity*, *symmetry*, and *transitivity*. It is interesting to note that for a partial ordering relation, reflexivity and transitivity are also assumed; it is only symmetry that gets changed to a diametrically opposed property which it is reasonable to call *anti-symmetry*.

An equivalence relation on A amounts to the same thing as a partitioning of A into disjoint sets. The details are set forth in Exercises 2 and 3.

The set of all ordered pairs (a, b) with $a \in A$ and $b \in B$ is called the Cartesian product of A and B and is written $A \times B$. The Cartesian product of n sets is likewise the set of ordered n -ples. To handle an infinite number of sets, we need to use an index set, as in Section 1.2. So let I be a set, and for each $i \in I$ let a set A_i be given. Then the Cartesian product $\prod_i A_i$ is the set of all "arrays" $\{a_i\}$ with $a_i \in A_i$. More formally, the Cartesian product is the set of all functions f from I to $\bigcup_i A_i$ which have the property $f(i) \in A_i$.

A relation between A and B can be described simply as a subset of $A \times B$.

EXERCISES

1. Discuss reflexivity, symmetry, and transitivity for the following relations: brother, sibling, friend, parent, ancestor, set-theoretic inclusion, greater than (for integers).
2. Suppose that on a set A there is given a partition: an expression of A as a union of disjoint subsets. For $a, b \in A$ define $a \sim b$ to mean that a and b lie in the same subset. Prove that this defines an equivalence relation.
3. Let \sim be an equivalence relation on A . For $a \in A$ define S_a to be the set of all b in A with $a \sim b$. Prove that two S_a 's are either disjoint or identical. Prove that the distinct S_a 's form a partition of A (they are called the *equivalence classes* of the relation).
4. On a set A let \leq be a relation which is a partial ordering except that anti-symmetry is not assumed. Define $a \sim b$ to mean that $a \leq b$ and $b \leq a$ both hold. Prove that \sim is an equivalence relation. Define a natural partial ordering on the equivalence classes. (As an application of this exercise, consider the relation of divisibility, defined on all integers, positive and negative.)
- 5.* Let f be a map of a set A into itself. Call a subset B of A *invariant* if $f(B) \subset B$. (More carefully, we ought to say invariant relative to f .) Call an invariant set *indecomposable* if it cannot be expressed as a union of disjoint nonvoid invariant sets. Prove that A has a unique expression as a disjoint union of indecomposable invariant sets. (*Hint*: Try to define an appropriate equivalence relation.)

2

Cardinal Numbers

2.1 COUNTABLE SETS

In this book we plan to treat finite numbers (the nonnegative integers) as known. Moreover, the distinction between finite and infinite sets will be treated intuitively; in fact, we have already done so several times. We shall concentrate on the novel idea that it may be possible, in a useful manner, to assign numbers to infinite sets.

Our starting point is the set N of positive integers, $N = \{1, 2, 3, \dots\}$. A moment's reflection is enough to convince us that N plays a crucial role in the study of infinite sets. Let A be any infinite set. Take any element of A and pair it with 1, take a second element and pair it with 2, etc. Since A is infinite, the process will not terminate in a finite number of steps. The end product of this procedure is a subset of A which has been placed in a one-to-one correspondence with N . (I can hear many readers snorting that, without any warning, I have sneaked in the countable axiom of choice. My reply is: guilty as charged. In an account of set theory designed for an apprentice mathematician, I think it out of place to fuss with the countable axiom of choice. Historical evidence is on my side. It was only when the close scrutiny of *uncountable* sets began that the axiom of choice got placed on the agenda.)

We shall say that a set is *countably infinite* if it can be put in a one-to-one correspondence with N , and we say that it is *countable* if it is either finite or countably infinite. Summarizing the remarks above, we note that *any infinite set contains a countably infinite subset*.

The next observation is that *any subset of a countable set is countable*. This is immediate: The crucial observation is that any infinite subset A of N can be placed in a one-to-one correspondence with N by counting up from the bottom, pairing the smallest element of A with 1, the second smallest with 2, etc.

We plan now to move from the positive integers to three successively larger sets: the integers, the rational numbers, and the algebraic numbers. The discussion is facilitated by a theorem (the first formal theorem of the book!).

THEOREM 1. *A countable union of countable sets is countable.*

Proof: We are given a countable index set I , and for each $i \in I$ a countable set A_i . We are to prove that $A = \bigcup_i A_i$ is countable. We shall suppose that I is infinite, noting at the appropriate point in the proof the slight change needed if I is finite. So I might as well be replaced by the set N of positive integers. The given countable sets are then designated by A_1, A_2, A_3, \dots

Since A_1 is countable, we can number off its elements, and we use the notation

$$a_{11}, a_{12}, a_{13}, \dots$$

There is the possibility that A_1 is finite. To handle this in a (more or less) uniform manner we put in harmless duplicates. Specifically, if A_1 has k elements, we list its elements as

$$a_{11}, a_{12}, \dots, a_{1k}$$

and agree that $a_{1,k+1}, a_{1,k+2}, \dots$ shall all be (say) a_{1k} . In the event that any later sets are finite, they are to be treated in the same way. With this agreement, we write

$$a_{21}, a_{22}, a_{23}, \dots$$

for the elements of A_2 , and continue in that fashion. The result will be the doubly infinite array

$$\begin{array}{cccc} a_{11} & a_{12} & a_{13} & \cdots \\ a_{21} & a_{22} & a_{23} & \cdots \\ \dots & & & \\ a_{n1} & a_{n2} & a_{n3} & \cdots \\ \dots & & & \end{array}$$

(If the index set is finite we terminate the array with the appropriate finite number of rows. The argument below is valid, the diagonals being suitably truncated; or the diagonals can be replaced by the columns. Alternatively, we can fill in all the rest of the rows with duplicates.) We now list the elements in a different order, using successive diagonals:

$$\begin{array}{ccccccc}
 a_{11} & a_{12} & a_{13} & a_{14} & \cdots & & \\
 \swarrow & & & & & & \\
 a_{21} & a_{22} & a_{23} & a_{24} & \cdots & & \\
 \swarrow & & & & & & \\
 a_{31} & a_{32} & a_{33} & a_{32} & \cdots & & \\
 \swarrow & & & & & & \\
 a_{41} & a_{42} & a_{43} & a_{44} & \cdots & &
 \end{array}$$

The list begins

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

There may be many duplicates in the list (in addition to the duplicates possibly introduced above, we did not assume the sets A_i to be disjoint). Keep only the first occurrence of an element, discarding all subsequent duplicates. If A is infinite, the final result exhibits a one-to-one correspondence between A and the positive integers. This concludes the proof of Theorem 1.

THEOREM 2. *The set of integers is countable.*

Proof: We exhibit the integers as the union of two countable sets:

$$\{0, 1, 2, 3, \dots\} \cup \{-1, -2, -3, \dots\}.$$

THEOREM 3. *The set of rational numbers is countable.*

Proof: Let Q be the set of rational numbers. For n a positive integer, let Q_n be the set of those rational numbers expressible with the denominator n . Thus Q_n is the set of all numbers m/n , m ranging over all integers. By Theorem 2, Q_n is countable. Since Q is the union of the sets Q_n , $n = 1, 2, 3, \dots$, it follows from Theorem 1 that Q is countable.

An *algebraic number* is a root of an equation with rational coefficients (it is equivalent to say integral coefficients). Thus $\sqrt{2}$, $\frac{1}{2}\sqrt[3]{3}$, and $\sqrt{2} + \sqrt{3}$ are all algebraic since they are roots of $x^2 - 2 = 0$, $8x^3 - 3 = 0$, and $x^4 - 10x^2 + 1 = 0$, respectively. Let it be officially agreed that we are talking about real algebraic numbers, although it would be harmless to admit complex numbers.

THEOREM 4. *The set of algebraic numbers is countable.*

Proof: Let T be the set of all algebraic numbers. Let S be the set of all polynomials with integral coefficients. For each $f \in S$, let $T(f)$ denote

the set of roots of f . Then each $T(f)$ is finite. Since T is the union of the sets $T(f)$, f ranging over S , it follows from Theorem 1 that it will suffice to prove that S is countable.

Each $f \in S$ has the form

$$f = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0,$$

with a_0, a_1, \dots, a_n integers. If $a_n \neq 0$ we say that f has degree n . Let S_n denote the set of polynomials of degree n . Since S is the union of the sets S_n , it follows from Theorem 1 that it suffices to prove that each S_n is countable. We fix one value of n .

Write

$$\|f\| = |a_n| + |a_{n-1}| + \cdots + |a_1| + |a_0|.$$

For k a positive integer, let S_{nk} denote the subset of S_n consisting of those f with $\|f\| \leq k$. Evidently S_n is the union of the sets S_{nk} , so (Theorem 1 again) it suffices to prove that each S_{nk} is countable. In fact, S_{nk} is obviously finite. Theorem 4 is proved.

At this point we might begin to think that all infinite sets are countable. This is not so; uncountable sets do exist. For instance, the set of all real numbers is not countable. To show this, we present the beautiful "diagonal argument" of Cantor.

Georg Cantor (1845–1913) was the founder of the theory of infinite numbers. He proved the uncountability of the real numbers in the first of a series of papers on set theory (*J. f. reine und angew. Math.* **77** (1874), 258–262). Actually, the proof of Theorem 5 given below is Cantor's second proof, as presented in a later paper (*Jahresbericht D. Math. Ver.* **1** (1892), 75–78). Cantor's first proof is sketched in Exercise 9. Although it is more abstract, it has the merit of not getting us involved in technicalities about decimal representations, and many readers might prefer it (I do). However, I do not venture to depart from the tradition of presenting Cantor's second proof at this point.

THEOREM 5. *The set of real numbers is not countable.*

Proof: It is slightly more convenient to prove the uncountability of the set S of real numbers between 0 and 1. Since any subset of a countable set is countable, this changed objective will suffice.

We think of the elements of S as infinite decimals. However, it is necessary to be honest about a small difficulty. The correspondence between S and infinite decimals is imperfect, since a decimal ending with an infinite string of 9's has a simpler representation. Thus,

$$0.429999 \dots \quad \text{and} \quad 0.430000 \dots$$

are the same real number. For definiteness, we agree not to use the version ending in 9's.

We turn to the proof that S is uncountable. Suppose on the contrary that S is countable, so that its members can be listed as a_1, a_2, a_3, \dots . We prove that this alleged enumeration of S must be incomplete by exhibiting a missing number b , and we make sure that b is missing by having b differ from a_1 in the first position after the decimal, differ from a_2 in the second position, etc. For a definite procedure let the n th digit of b be 2 if the n th digit of a_n is 1, and let the n th digit of b be 1 otherwise. (Any such procedure will do, but to avoid any question of ambiguity stay away from 9's and 0's.) We give an illustration. Suppose that

$$\begin{aligned} a_1 &= 0.4275 \dots, \\ a_2 &= 0.0193 \dots, \\ a_3 &= 0.9912 \dots, \\ a_4 &= 0.1234 \dots, \end{aligned}$$

Then $b = 0.1221 \dots$. Certainly b is missing from the list, since it differs from each a_i in at least one of its digits. This contradiction proves Theorem 5.

Remark: A naive thought that might come up is the following: Although we missed b in this effort, we could change our minds and place b at the top of the list. Perhaps the list b, a_1, a_2, \dots , exhausts S . But of course the argument will apply again to produce a fresh element c differing from b and all the a 's.

A real number is called *transcendental* if it is not algebraic. Since the set of real numbers is uncountable, while the set of algebraic numbers is countable, it follows that transcendental real numbers must actually exist. In fact, the set of transcendental real numbers is clearly uncountable, so that we have proved that "hordes" of transcendentals exist. It is often said that Cantor's proof is not "constructive," and so does not yield a tangible transcendental number. This remark is not justified. If we set up a definite listing of all algebraic numbers (it is feasible to do this by examining the proof of Theorem 4), and then apply the diagonal procedure of Theorem 5, we get a perfectly definite transcendental number (it could be computed to any number of decimal places). The method of Exercise 9 could be used in a similar way. (I owe these remarks to R. M. Robinson.) But this does not diminish our interest in knowing the status of "concrete" numbers like e and π . To this day it is a fairly tricky matter to prove that e and π are transcendental. (The transcendence of e was proved by Hermite in 1873; that of π was proved by Lindemann in 1882.) Hilbert's seventh problem was to prove that α^β is transcendental if α and β are algebraic with $\alpha \neq 0, 1$ and β irrational. (Examples: $2^{\sqrt{2}}$, $i^i = e^{-\pi/2}$.) This was proved by Gelfond in 1934, and independently by Schneider shortly thereafter. Nobody knows the status of π^π . Perhaps it is even rational.

Charles Hermite (1822–1901) made numerous important contributions in many branches of mathematics. Ferdinand Lindemann (1852–1939), in addition to proving π transcendental, ranks as important since he was Hilbert's teacher.

David Hilbert (1862–1943) was, by unanimous agreement, the greatest mathematician of his time. At the 1900 Paris International Mathematical Congress he propounded 23 problems which he considered important for the future of mathematics. (An English translation appears in the *Bull. Amer. Math. Soc.* 8 (1902), 437–479.)

Aleksandr Osipovich Gelfond (1906–1968) was a leading authority on transcendental numbers. A book on Hilbert's problems was published in 1969 in the U.S.S.R.; its seventh chapter was written by Gelfond and published posthumously.

EXERCISES

1. Let A be a countable set and suppose there is given a function mapping A onto a set B . Prove that B is countable.
2. Prove that the Cartesian product of two countable sets is countable. Generalize to the Cartesian product of a finite number of sets.
3. Prove that the set of all finite subsets of a countable set is countable. (*Hint*: Use Exercises 1 and 2.)
4. Show that the set N of positive integers can be represented as a union $N = \cup A_i$ of an infinite number of disjoint infinite sets.
- 5.* (a) A subset M of a partially ordered set L is said to be *cofinal* in L if for any $x \in L$ there exists $y \in M$ with $x \leq y$. Prove that any countable lattice has a cofinal chain.
(b) Let A be an uncountable set, and let L be the partially ordered set consisting of all finite subsets of A (the ordering on L is set-theoretic inclusion). Prove that L does not have a cofinal chain.
- 6.* A chain C is said to be dense in itself if the following holds: For any a, b in C with $a < b$ there exists $x \in C$ with $a < x < b$. Let A and B be countable chains each of which is dense-in-itself and furthermore does not have a top or bottom element. Prove that A and B are order-isomorphic. (*Hint*: Build an order-isomorphism between A and B step by step. To make sure that A and B are entirely accounted for, two precautions are necessary. The elements of A and B should be numbered off in advance (this numbering has nothing to do with the given orderings on A and B). Secondly, the procedure should *alternate* between A and B .)
- 7.* Let A be a countable chain which is dense-in-itself. Prove that, up to order-isomorphism, there are exactly four possibilities for A . (*Hint*: Modify Exercise 6 appropriately.)

- 8.* Prove that any countable chain is order-isomorphic to a chain of rational numbers. (*Hint*: Build the order-isomorphism stepwise. In contrast to Exercises 6 and 7, no alternation is needed.)
- 9.* Let L be a conditionally complete dense-in-itself chain. (See Exercise 11 in Section 1.3 for the definition of “conditionally complete.”) Prove that L is uncountable. (*Hint*: Suppose the contrary and fix an enumeration of L . Starting with the first element x_1 (which we can suppose not to be a maximal element of L), let y_1 be the first element with $y_1 > x_1$, x_2 the first subsequent element with $x_1 < x_2 < y_1$, y_2 the first subsequent element with $x_2 < y_2 < y_1$, etc. In this way an ascending sequence of x 's and a descending sequence of y 's are constructed, looking like

$$x_1 < x_2 < \cdots < x_n < \cdots < y_n < \cdots < y_2 < y_1.$$

If z is the least upper bound of the x 's, the element z must have been missed.

As noted above, this is Cantor's original proof of the uncountability of the real numbers. *Remark*: We only need least upper bounds for *countable* sets bounded from above.)

10. Let A be an infinite set, B a finite subset of A , and C the complement of B in A . Prove that there exists a one-to-one correspondence between A and C .
- 11.* Let A be an uncountable set, B a countable subset of A , and C the complement of B in A . Prove that there exists a one-to-one correspondence between A and C .
- 12.* Let L be the partially ordered set of all subsets of a countably infinite set (ordered by inclusion). Prove that L contains a chain order-isomorphic to the chain of all real numbers. (*Hint*: Number off the rationals and assign elements of L to the rationals stepwise, always allowing an “infinite amount of room” between two choices. Take least upper bounds to insert the real numbers.)
- 13.* Let A be a countably infinite set. Show that A contains continuum many subsets (i.e., a collection of subsets in one-to-one correspondence with the real numbers) such that any two have a finite intersection. (*Hint*: Think of A as the rational numbers. For each irrational number, pick a sequence of rationals converging to it.)
14. Let L be a partially ordered set in which every countable chain has an upper bound. Let S be a countable subset with the following property: For any a and b in S there exists c in S with $a \leq c$, $b \leq c$. Prove that S has an upper bound in L .

2.2 CARDINAL NUMBERS

We proceed to come to grips with the project of assigning numbers to infinite sets. Actually we shall use the designation “cardinal number,” in

order to distinguish these numbers from the "ordinal numbers" which will follow in due course.

Whatever the concept of cardinal number finally turns out to be, we firmly assert that it is to be an object attached to sets in such a way that two sets have the same cardinal number if and only if there is a one-to-one correspondence between them. Obvious though this seems today, it was very novel and thought-provoking when Cantor proposed it. Cantor's words are worth quoting verbatim (*Math. Annalen* 46 (1895), 481).

Unter einer "Menge" verstehen wir jede Zusammenfassung M von bestimmten wohlunterschiedenen Objekten m unserer Anschauung oder unseres Denkens (welche die 'Elemente' von M genannt werden) zu einer Ganzen.

. . .

"Mächtigkeit" oder "Cardinalzahl" von M nennen wir den Allgemeinbegriff, welcher mit Hilfe unseres activen Denkvermögens dadurch aus der Menge M hervorgeht, dass von der Beschaffenheit ihrer verschiedenen Elemente m und von der Ordnung ihres Gegehenseins abstrahirt wird.

The German language seems particularly suitable for discourse of this kind. An English translation by Jourdain is available in *Contributions to the Founding of the Theory of Transfinite Numbers*, published in 1915 and reprinted by Dover in 1952.

Here is a plausible way to organize our thoughts on the subject of one-to-one correspondences and cardinal numbers. Write $A \sim B$ if there exists a one-to-one correspondence between the sets A and B . Evidently this relation is reflexive, symmetric, and transitive, and so it is an equivalence relation. As such it partitions all sets into equivalence classes. Call these equivalence classes cardinal numbers. The cardinal number of a set is then to be the equivalence class to which it belongs. We can take this seriously for finite sets as well as infinite sets, thus learning, for instance what 2 is: The number 2 is the equivalence class consisting of all pairs.

There are two conceivable objections. The definition makes cardinal numbers highly intangible objects, and there is some danger of flirting with the paradoxes to be discussed quite soon.

Let us instead take an informal *ad hoc* attitude, which we shall find serves our purposes quite well. We shall not say what a cardinal number actually is. The cardinal number of a set A is a "thing" which we write $\alpha(A)$. We have $\alpha(A) = \alpha(B)$ if and only if there exists a one-to-one correspondence between A and B .

The first infinite cardinal on the agenda is obviously $\alpha(N)$, where N is the set of positive integers. The notation for $\alpha(N)$ is \aleph_0 , the symbol \aleph being the first letter of the Hebrew alphabet. A set A satisfies $\alpha(A) = \aleph_0$ if and only if A is countably infinite.

We saw in the last section that the set of real numbers has a cardinal number different from \aleph_0 . So we have a second cardinal number; we write c for it (c stands for "continuum," and c is often called the cardinal number of the continuum).

We have two infinite cardinal numbers. Are there any more? The crucial point in getting an affirmative answer is another theorem of Cantor, showing that for any set A the power set $P(A)$ has a "larger" cardinal number than that of A (we shall give a precise meaning to the adjective "larger" quite soon). The proof of Theorem 6 is in essence the same diagonal procedure as was used in proving Theorem 5. But any striking proof is worth seeing twice.

We prefer to state Theorem 6 directly for sets, without reference to their cardinal numbers.

THEOREM 6. *For any set A , there does not exist a function mapping A onto its power set $P(A)$.*

Remarks: 1. *A fortiori*, there cannot exist a one-to-one correspondence between A and $P(A)$. When we have assembled all the facts concerning cardinal numbers, it will be plain that the statement in Theorem 6 is not really stronger. But it is pleasant at this point to collect the maximum amount of information from the proof.

2. Theorem 6 is really interesting only for infinite sets. But it is true for finite sets, and the proof to be given is valid. To be sure the inequality $n < 2^n$ for nonnegative integers is not very exciting.

Proof: We make an indirect proof, and therefore assume the existence of a function $f: A \rightarrow P(A)$ which is onto. Thus to any element a in A there is associated $f(a)$, a subset of A . Now we distinguish two possibilities: a may or may not be a member of its corresponding $f(a)$. Let us collect in a subset B of A all a such that $a \notin f(a)$. (*Query:* How does the rest of the proof work out if instead we try the set of all a with $a \in f(a)$?)

We interpolate an illustration. Suppose that $A = \{1, 2, 3\}$. Then the elements of A might, for instance, be sent to subsets of A as follows:

$$\begin{aligned} 1 &\rightarrow \{2\} \\ 2 &\rightarrow \{2, 3\} \\ 3 &\rightarrow \{1, 2\}. \end{aligned}$$

(Of course, since Theorem 6 is true, we cannot exhibit the three images filling up the eight-member power set of A !) The elements which are not members of their corresponding sets are 1 and 3; thus $B = \{1, 3\}$. Note that $\{1, 3\}$ does not come from any element of A .

We return to the proof of Theorem 6. Like any subset of A , B must come from such element of A . We write b for such an element, so that $f(b) = B$. Can this element b lie in B ? If so, by the very definition of B , we have $b \notin f(b)$, which is nonsense. So it must be the case that b is not in B . But again we are in trouble, for by the definition of B , $b \in f(b)$ follows, again a contradiction. Thus we must admit that the function f cannot exist.

Several years after Cantor gave the intriguing proof of Theorem 6, Bertrand Russell asked the following question: What if A is already the universal set, i.e. the set of which absolutely everything is a member? How can $P(A)$ be "larger" than A in that case?

Bertrand Russell (1872–1970) was one of the best known of men during his lifetime. His early work on mathematical logic was climaxed by the monumental *Principia Mathematica*, written with Alfred North Whitehead. For the paradox, see page 101 of his *Principles of Mathematics*. He died while a draft of this manuscript was being typed.

The result was the celebrated Russell paradox. To present it, we begin by admitting that it is conceivable that a set A might be a member of itself: $A \in A$. To be sure, the possibility is a little bizarre. So it is reasonable to call a set "pleasant" if this is not the case. In other words, A is pleasant if $A \notin A$. Now form R , the set of all pleasant sets. Either $R \in R$ or $R \notin R$ must be the case. But, just as in the proof of Theorem 6, either assumption traps us in a contradiction.

Numerous variants of the Russell paradox have been invented. Here is a verbal one that nonmathematicians can enjoy. Consider adjectives in the English language, and ask whether or not they are self-descriptive. Example: The word "short" is short (or reasonably so), but the word "long" is certainly not long. So "short" is self-descriptive and "long" is not. Another example: The word "English" is English, but the word "German" is not German (it is English). Call an adjective *homological* if it is self-descriptive, otherwise *heterological*. Surely any adjective is one or the other. Now try out the adjective "heterological".

(I recall learning this version of the paradox in a course given by Quine at Harvard. It is referred to on page 254 of his *Set Theory and Its Logic* (Harvard, 1969) where it is credited to Grelling.)

It is amusing to write the following on one side of a card: The statement on the other side of this card is true. Then write the same thing on the other side, except that the word "true" is replaced by "false." What is the truth status of the two statements?

An absurdly simple version of the paradox occurs when a speaker says, "The sentence I am now uttering is a falsehood."

Russell's paradox shows that if set theory is treated too naively it can lead to disaster. It is beyond the scope of this book to take up careful axioms for set theory, but we shall mention the method commonly used to circumvent paradoxes. Based on an idea of von Neumann, it distinguishes "good" sets and "bad" sets. Good sets are allowed to be members of other sets, but bad sets are blackballed from membership. The Russell set R of all nonself-members can be legally formed. The argument for the paradox does not lead us into irrevocable trouble, but merely shows that R is a bad set.

John von Neumann (1903–1957) was one of the great mathematicians of the twentieth century. He contributed significantly to every field he touched (logic, Hilbert space, game theory, computers, etc.). The paper referred to appeared in the *J. f. reine und angew. Math.* **154** (1925), 219–240.

EXERCISES

1. Prove that the set of positive real numbers has cardinal number c . (*Hint:* Try the map $x \rightarrow e^x$.)
2. Prove that the open unit interval (the set of all real x with $0 < x < 1$) has cardinal number c . (*Hint:* Use Exercise 1 and the map $x \rightarrow x/(1+x)$, defined on all positive real numbers.)
3. What is the cardinal number of the closed unit interval (all x with $0 \leq x \leq 1$)? Of the half open interval (all x with $0 \leq x < 1$)?
4. Extend Exercises 2 and 3 to the interval from a to b , where a and b are real numbers with $a < b$.
5. What is the cardinal number of the set of irrational real numbers? Of the set of transcendental real numbers?

2.3 COMPARISON OF CARDINAL NUMBERS; ZORN'S LEMMA

We turn to a more systematic examination of the comparison of cardinal numbers. We shall typically use d and e (decorated with subscripts as necessary) to denote cardinal numbers.

DEFINITION. Let d and e be cardinal numbers. Let D and E be sets with $\alpha(D) = d$, $\alpha(E) = e$. We say that $d \leq e$ if there exists a one-to-one function from D into E . (This is obviously independent of the choice of the representatives D and E .) We say that $d < e$ if $d \leq e$ and $d \neq e$.

As illustrations we observe (i) $n < d$ if n is finite and d is infinite, (ii) $\aleph_0 \leq d$ for any infinite d , (iii) $\aleph_0 < d$ for any infinite uncountable d .

Warning: Suppose there exists a one-to-one function from D onto a proper part of E . We cannot conclude that $d < e$, since $d = e$ is perfectly possible (this is in fact characteristic of infinite sets). So to convince ourselves that $d < e$ we have to argue that there is a one-to-one function from D into E and that furthermore there is no possible one-to-one function from D onto E .

The relation \leq on cardinal numbers satisfies reflexivity and transitivity. In other words we have $d \leq d$ for all d , and if $d_1 \leq d_2$, $d_2 \leq d_3$ then $d_1 \leq d_3$. These statements are very easy to check, and we leave their verification to the reader. There remains only the property of antisymmetry in order to see that cardinal numbers form a partially ordered set. We delay this for a moment.

Let D be a set with cardinal number d . We write 2^d for the cardinal number of the power set $P(D)$. This is a sensible concept, for it is obvious that if two sets have the same cardinal number then their power sets have the same cardinal number. In Section 2.6 we shall explain why the notation 2^d is appropriate, and we shall also prove that $2^{\aleph_0} = c$.

There is an evident one-to-one map from any set A into its power set $P(A)$: Just send $a \in A$ into the set $\{a\}$. Hence Theorem 6 gives us $d < 2^d$ for any cardinal number. Applying this in turn to the cardinal number 2^d , we furthermore have $2^d < 2^{2^d}$. From the two statements $d < 2^d$ and $2^d < 2^{2^d}$ can we conclude that $d < 2^{2^d}$? We must not be hasty here. Indeed a moment's reflection shows that proving $<$ to be transitive is exactly the same as proving \leq to be antisymmetric. This is true for cardinal numbers, but it is a substantial theorem, recorded subsequently as Theorem 9 (Theorems 7 and 8 are equivalent to Theorem 9). However, in the case at hand we have extra information that resolves the difficulty. We proved in Theorem 6 that for no set A does there exist a map of A onto its power set (whether or not the map is one-to-one). From this we can deduce not just that 2^d is unequal to d , but also that $2^d < d$ is ruled out. Hence we can in fact exclude $d = 2^{2^d}$, since it would imply $2^d < 2^{2^d} = d$. Thus we have $d < 2^{2^d}$. Continuing in this way we can produce the infinite sequence of cardinal numbers

$$d, 2^d, 2^{2^d}, 2^{2^{2^d}}, \dots$$

with the property that each is less than any of the cardinal numbers to the right of it.

It is a charming aspect of set theory that antisymmetry of \leq for cardinal numbers admits an elementary proof that can be given very early in the

subject. The theorem was missed by Cantor (who was very anxious to prove it). It was proved independently by Bernstein and Schröder. See pages 449–450 of Cantor's *Collected Works* for an interesting exchange of correspondence concerning it between Cantor and Dedekind.

The theorem that needs to be proved may be stated as follows:

THEOREM 7. *Let A and B be sets such that there exists a one-to-one map of A into B and a one-to-one map of B into A . Then there exists a one-to-one correspondence between A and B .*

An unusually clear and readable proof of Theorem 7 appears on page 340 of the third edition of *A Survey of Modern Algebra* by Birkhoff and Mac Lane. For a reader who would like a brisk but rather sophisticated proof, Exercise 3 is recommended (note that Exercise 15 in Section 1.4 must be done as a prelude). In this account still another proof is offered. It is not fundamentally different from the usual proofs; indeed it is unlikely that anyone will ever devise a truly novel proof. Nor can the proof be recommended for its brevity. It does however have the merit of being a rather quick corollary of material that is worth developing for its own sake. Once it has been decided to work up this theory of infinite cycle decompositions, the procedure flows naturally and needs no artifices.

We shall deduce Theorem 7 from another theorem which is essentially just a variant.

THEOREM 8. *Let f be a one-to-one map of a set A into itself. Let C be a subset of A containing $f(A)$. Then there exists a one-to-one correspondence between A and C .*

We give at once the deduction of Theorem 7 from Theorem 8.

Proof of Theorem 7 from Theorem 8: Let $r: A \rightarrow B$ and $s: B \rightarrow A$ be the given functions, and write $f = sr$. Then f is a one-to-one function of A into itself. Let $C = s(B)$. Then $C \supset f(A)$ and we are in a position to apply Theorem 8, getting a one-to-one correspondence between A and C . Since s provides a one-to-one correspondence between B and C , we achieve the desired one-to-one correspondence between A and B .

Very likely, nearly all readers are familiar with the decomposition of a permutation of a finite set into disjoint cycles. We shall review it briefly.

Let f be a permutation of a finite set A , so that f maps A one-to-one onto itself. Elements a_1, \dots, a_r in A form a *cycle* under f if $f(a_1) = a_2$, $f(a_2) = a_3, \dots, f(a_{r-1}) = f(a_r), f(a_r) = a_1$. In words: f sends each of a_1, \dots, a_r into the next and sends a_r back into a_1 . The set A splits into

disjoint cycles. To see this, start with an arbitrary a in A and apply f to it repeatedly. The sequence obtained returns to a after a finite number of steps and we have thus constructed a cycle. If this cycle does not exhaust A , start with a new element and treat it in the same fashion. The procedure is continued till the disjoint cycles obtained fill up A .

What we now wish to do is to repeat this discussion with A allowed to be infinite; then we shall broaden the context a trifle further by not insisting that f is onto (while maintaining the assumption that it is one-to-one). This calls for the introduction of infinite cycles.

Our notation for an infinite cycle will be

$$(\dots a_{-2} a_{-1} a_0 a_1 a_2 \dots)$$

and the understanding is that the function f sends every element into the next one on the right, so that we have $f(a_i) = a_{i+1}$ for all integers i (positive, zero, and negative). More exactly, this will be called a *bilateral* infinite cycle, to be distinguished from the *unilateral* ones that will shortly be introduced.

Now the following is true in the infinite case just as in the finite case: If f is a one-to-one mapping of a set A onto itself, then A splits under f into disjoint cycles. The method of obtaining the decomposition follows the same lines in the infinite case as in the finite case. Start with any $a \in A$, and apply to it repeatedly both f and f^{-1} . The array that is generated can be exhibited as

$$\dots, f^{-1}(f^{-1}(a)), f^{-1}(a), a, f(a), f(f(a)), \dots$$

If there is ever a repetition in this array, the whole collection of elements boils down to a finite cycle containing a . Otherwise we obtain a bilateral infinite cycle. By repeating this procedure, we insert every element into a cycle, and different cycles are disjoint.

We hasten to reassure any worried reader that no transfinite procedure is called for here. The decomposition can be done all at once, and in fact the discussion serves as a nice example of an equivalence relation. Introduce the notation f^n to mean the result of composing f with itself n times, and extend it to all integers by having f^0 mean the identity function and $f^{-m} = (f^{-1})^m$. Then the equivalence relation we need is definable as follows: $a \sim b$ if $b = f^n(a)$ for some n . The equivalence classes are precisely the cycles discussed above.

We proceed to the final step. By a *unilateral* infinite cycle

$$(a_1, a_2, a_3, \dots)$$

we mean a subset where $f(a_i) = a_{i+1}$ for $i = 1, 2, 3, \dots$ and a_1 is not in the range of f . Now suppose that f is a one-to-one mapping of A into itself that need not be onto. We assert that A splits into disjoint cycles, where the concept has been widened to include unilateral infinite cycles.

The discussion is just a slight variant of what we have already done. Starting with an element $a \in A$, we apply f to it repeatedly. Although f is not onto, we shall (in this discussion) venture to use the symbol f^{-1} ; it is defined only on the range of f . Apply f^{-1} to a as long as possible. This may go on forever, or may end in a finite number of steps. If we ever encounter a repetition, a finite cycle containing a will emerge. Otherwise we get an infinite cycle, which may be either bilateral or unilateral. The decomposition is again describable by an equivalence relation.

With all this accomplished we are ready for the proof of Theorem 8.

Proof of Theorem 8: Since Theorem 8 was stated several pages back, we recapitulate it. We are given a set A , a one-to-one map f of A into itself, and a set C lying between A and $f(A)$. We are to devise a one-to-one map of A onto C .

Break A into cycles under f . We shall use the following notation: $A = D \cup E$, where D combines all the finite cycles and bilateral infinite cycles, and E combines all the unilateral infinite cycles. Note that f maps D onto itself. (In fact, D is the largest subset of A carried onto itself by f . See in this connection Exercise 17 in Section 1.4.) In order to discuss E , let us list what might be some of the unilateral infinite cycles:

$$\begin{aligned} &(a_1, a_2, a_3, \dots) \\ &(b_1, b_2, b_3, \dots) \\ &(c_1, c_2, c_3, \dots) \end{aligned}$$

We can now envisage $f(A)$ very explicitly, for $f(A)$ consists of D together with all of E except the initial elements of the unilateral cycles. In symbols.

$$f(A) = D \cup \{a_2, a_3, \dots\} \cup \{b_2, b_3, \dots\} \cup \{c_2, c_3, \dots\} \cup \dots$$

while a_1, b_1, c_1, \dots do not belong to $f(A)$.

What can C look like? The answer is that C adds to $f(A)$ some of the missing initial elements of unilateral cycles. For instance, C might adjoin to A the elements a_1 and c_1 , but might not contain b_1 . We can now invent the required one-to-one correspondence (say g) between A and C . On D take g to be any one-to-one map of D onto itself; for instance $g = f$ or $g =$ the identity will do. On the unilateral cycles which appear completely in C , take g to be the identity. Finally, on the unilateral cycles which remain incomplete in C , take $g = f$. Then g maps A one-to-one onto C . This proves Theorem 8, and thereby also Theorem 7.

Remark: The decomposition into cycles obtained here is a special case of Exercise 5 in Section 1.5.

We restate Theorem 7.

THEOREM 9. *The relation \leq on cardinal numbers is a partial ordering.*

We immediately ask: Is the partial ordering of cardinal numbers a linear ordering? In other words, can any two sets be compared?

Let us launch a simple-minded attack on the problem. Someone hands us two infinite sets, A and B . We proceed to pair off the members of A and B . An arbitrarily chosen element of A is paired with an arbitrary element of B . These are put aside, and second elements of A and B get paired. The process is continued forever. At this point we note that if A is countable, the selections could have been made so as to exhaust A . So further activity is really needed only when A and B are both uncountable. We start up the pairing off procedure again, and let it run forever again. Stretching our minds further, we envisage the business of going on forever repeated a "forever number of times."

Fatigue is setting in and we are tempted to wave a hand and announce that "ultimately" either A or B must get used up. Probably no mathematician anywhere, any time, was willing to let it go at that. The instructions to go on forever, do this forever times, do this in turn forever times, . . . and the claim that this will somehow end, seem to leave everyone feeling very uncomfortable.

Here is a capsule history of the development of the subject. Thoughts like the preceding led Cantor to systematize the business of going on forever repeatedly, by introducing the concept of well-ordering. (In Chapter 3 we shall discuss well-ordering and the axiom of choice in some detail; at present the reader should regard this paragraph as gossip that can be omitted.) The crucial question then became: Can every set be well-ordered? Zermelo supplied a proof, and had the genius to see that a special axiom which he called the axiom of choice played a critical role. Years later an alternate technique, often simpler to use, became standard. In this form the key axiom is popularly called Zorn's lemma.

Max Zorn began his career in the field of alternative rings. The paper of his that is relevant to the present discussion appeared in the *Bull. Amer. Math. Soc.* 41 (1935), 667–670. It was influential in making the new method well known. One of the forms of Zorn's lemma appeared in 1914 on page 140 of the first edition of Hausdorff's *Mengenlehre*, a basic treatise on set theory which remains valuable to this day.

Felix Hausdorff (1868–1942) began his career in astronomy and subsequently made important contributions in analysis as well as in set theory. Faced with deportation by the Nazis, he and his wife committed suicide on January 26, 1942 in Bonn, where he was a professor emeritus.

Ernst Zermelo (1871–1953) followed up his proof that the axiom of choice implies well-ordering with basic studies on axiomatic set theory which marked a turning point in mathematical logic. He edited Cantor's collected

works, published in 1932. For refusal to give the Nazi salute and insults to the Führer he was threatened with dismissal from his post at the University of Freiburg. In reply he resigned on March 2, 1935. In 1946 he was reappointed.

The name "Zorn's lemma" is misleading. A lemma is a little theorem, usually proved on the way to a big theorem. But, little theorem or big theorem, the recommended attitude to Zorn's lemma is that it is an *axiom*. See Section 3.3 for more on this.

What we shall now do is set down one convenient form of Zorn's lemma and proceed forthwith to use it in the style that is customary today.

ZORN'S LEMMA: *Let L be a partially ordered set in which every chain has an upper bound. Then L contains a maximal element.*

Remarks: 1. The upper bound in question need not be in the chain, but of course it must be in L .

2. By a maximal element x we mean one such that $x < y$ is not true for any y in L . x need not top all the elements of L . As an extreme case, L might be totally unordered, in which case every element of L is maximal.

3. Let us get an intuitive idea of why this is a plausible axiom. Start with any $a_1 \in L$. If a_1 is maximal our search is over. Otherwise we have $a_2 \in L$ with $a_1 < a_2$. If a_2 is not maximal, we have a larger a_3 . If this goes on forever we construct a chain

$$a_1 < a_2 < a_3 < \cdots < a_n < \cdots$$

which by hypothesis has an upper bound. This upper bound can be used to keep the process going. To be sure, the question of how or why the procedure is going to end is just as puzzling as it was in our discussion of the comparability of two sets. But it is less troubling to be puzzled by an axiom than by a proof. In Section 3.3 our faith in the axiom will get some reinforcement.

We proceed to the proof of comparability via Zorn's lemma. In this first proof we shall give very complete details. Later proofs will be appropriately streamlined.

A reader familiar with a little abstract algebra may find it a good idea to try out some of the examples in Appendix 2 as his first illustrations of Zorn's lemma. Exercise 4 is also recommended. For the logical development of the subject it is presumably a good idea to prove Theorem 10 at this point. However, it is undeniable that the length of the proof is a little discouraging.

THEOREM 10. *For any sets A and B there exists either a one-to-one function of A into B or a one-to-one function of B into A (or both).*

Proof: We offer a preliminary piece of advice concerning the strategy of most proofs that use Zorn's lemma. Normally, the partially ordered set is taken to be the set of all "intermediate stages" appropriate to the problem at hand. In the present case the intermediate stages can conveniently be taken to be triples (A_i, B_i, f_i) , where A_i is a subset of A , B_i is a subset of B , and f_i is a one-to-one mapping of A_i onto B_i . Thus our partially ordered set L is the set of all such triples. To define the partial ordering on L , let a second triple (A_j, B_j, f_j) be given. We say that $(A_i, B_i, f_i) \leq (A_j, B_j, f_j)$ if $A_i \subset A_j$, $B_i \subset B_j$, and f_j restricted to A_i coincides with f_i . We argue that this defines a partial ordering on L . Reflexivity and antisymmetry are immediate. Transitivity is also easy, but we give a little detail. The crucial point is this: We are given $A_i \subset A_j \subset A_k$ and f_i, f_j, f_k defined on these three sets. We know that f_j restricted to A_i equals f_i , and that f_k restricted to A_j equals f_j . We have to verify that f_k restricted to A_i equals f_i . Take $x \in A_i$. Then $f_j(x) = f_i(x)$. Since $x \in A_j$ we also have $f_k(x) = f_j(x)$. Hence $f_k(x) = f_i(x)$, as required.

Now we need to check that any chain C in L has an upper bound. We get the desired upper bound in the obvious way. Let the elements of C be the triples $(A_\alpha, B_\alpha, f_\alpha)$, α running over some index set. Take $A_0 = \cup A_\alpha$, $B_0 = \cup B_\alpha$. Our problem is to define a suitable function f_0 on A_0 . If x is an element of A_0 , we observe that x got into A_0 by virtue of being in some A_α . There the definition $f_0(x) = f_\alpha(x)$ awaits us. The urgent thing is to check that this definition does not depend on the choice of α . So suppose that x also lies in A_β . Because C is a chain we have inclusion one way or the other between the triples $(A_\alpha, B_\alpha, f_\alpha)$ and $(A_\beta, B_\beta, f_\beta)$. Say for definiteness that $(A_\alpha, B_\alpha, f_\alpha) \leq (A_\beta, B_\beta, f_\beta)$. Then $A_\alpha \subset A_\beta$ and f_β restricted to A_α coincides with f_α . This tells us that $f_\beta(x) = f_\alpha(x)$, as we needed. In this way we have a well defined function f_0 on A_0 . There are now several things to be verified.

f_0 is one-to-one: Suppose that $f_0(x) = f_0(y)$ for x and y in A_0 . We have, say, $x \in A_\gamma$ and $y \in A_\delta$. As before, A_γ and A_δ are comparable. So we may assume that both x and y lie in A_γ . Then f_0 on x and y is given by f_γ , and f_γ is known to be one-to-one on A_γ .

f_0 maps A_0 into B_0 : Given $x \in A_0$, we have that x lies in some A_α . Then $f_0(x) = f_\alpha(x)$ and $f_\alpha(x) \in B_\alpha \subset B_0$.

f_0 maps A_0 onto B_0 : Take $z \in B_0$. We have that z lies in some B_α . Since f_α maps A_α onto B_α , we have $w \in A_\alpha$ with $f_\alpha(w) = z$. Then $w \in A_0$ and $f_0(w) = z$.

For any α , f_0 restricted to A_α is f_α : We take $x \in A_\alpha$ and inquire whether $f_0(x) = f_\alpha(x)$. This is so by the definition of f_0 .

All these facts add up to the statement that the triple (A_0, B_0, f_0) satisfies $(A_\alpha, B_\alpha, f_\alpha) \leq (A_0, B_0, f_0)$ for every α . We have thereby constructed the requisite upper bound for the chain C .

So the partially ordered set L fulfills the hypothesis of Zorn's lemma. Therefore L has a maximal element. Call it (A^*, B^*, f^*) . We claim that either A^* is all of A or B^* is all of B . For suppose that both statements are false. Then we can select u in A but not in A^* , and v in B but not in B^* . Now define g on $\{A^*, u\}$, the set consisting of the members of A augmented by u , by making it coincide with f^* on A^* and setting $g(u) = v$. Then the triple

$$(\{A^*, u\}, \{B^*, v\}, g)$$

is properly larger than (A^*, B^*, f^*) , contradicting maximality. The proof of Theorem 10 is thereby concluded.

Once the general plan of attack of such a proof is laid out, the details are largely routine. Therefore (as we remarked above) future applications of Zorn's lemma will be presented in a brisker style.

We restate Theorem 10.

THEOREM 11. *Cardinal numbers form a chain under \leq .*

EXERCISES

- In each of the following cases a one-to-one mapping of a set into itself is given. Describe the cycle structure.
 - $x \rightarrow x + 1$ on the positive integers.
 - $x \rightarrow x + 2$ on the positive integers.
 - $x \rightarrow x + n$ on the positive integers, n being a positive integer.
 - $x \rightarrow x + n$ on all integers, n being any integer.
 - $x \rightarrow 2x$ on the positive integers.
 - $x \rightarrow -x + 2$ on all integers.
 - $x \rightarrow x^{-1}$ on all nonzero rational numbers.
- Let L be a lattice in which every chain has an upper bound. Prove that L has a unique maximal element.
- * (a) Let functions $f: A \rightarrow B$ and $g: B \rightarrow A$ be given. Define a mapping $h: P(A) \rightarrow P(A)$ by $h(S) = (g[f(S)'])'$ for S a subset of A (the inner prime is complementation with respect to B ; the outer is complementation with respect to A). Prove that $S \supset T$ implies $h(S) \supset h(T)$.
 (b) Use part (a) and Exercise 15 in Section 1.4 to give another proof of Theorem 7.

4. Prove that any partially ordered set contains a maximal totally unordered subset. (*Explanation:* Let L be the partially ordered set. A subset A of L is totally unordered if $x \leq y$ for x, y in A implies $x = y$. A is maximal, relative to this property, if no subset of L properly containing A is totally unordered.)
- 5.* Prove that the order on an arbitrary partially ordered set can be strengthened so as to make it a chain. (Strengthening a partial order means keeping all instances of \leq and perhaps adding more. Make a partially ordered set out of all ways of strengthening the order. Argue that Zorn's lemma is applicable. That the resulting maximal object must be a chain is seen by proving that if a and b are incomparable then $a \leq b$ can be introduced, contradicting maximality.)
- 6.* Let L be a lattice in which every chain has a least upper bound and a greatest lower bound. Prove that L is complete.

2.4 CARDINAL ADDITION

Let d and e be cardinal numbers. To define $d + e$ we take disjoint sets D and E with $o(D) = d$, $o(E) = e$, and set $d + e = o(D \cup E)$. The definition is clearly independent of the choice of representatives D and E . Note that it gives the right answer in the finite case.

One small objection might be raised. Can we be certain that disjoint representatives D and E can be found? A rather ridiculous device can quiet our fears on this score. First take any representatives D and E . Then replace D by the set D^* of ordered pairs $(x, 1)$, x ranging over D . Likewise replace E by the set E^* of ordered pairs $(y, 2)$, y ranging over E . Then D^* and E^* are disjoint, and they can serve as representatives of d and e .

It turns out (Theorem 14) that addition of infinite cardinal numbers gives nothing new. We first prove two preliminary theorems.

THEOREM 12. *Any infinite set can be expressed as a disjoint union of countably infinite subsets.*

Proof: Let A be the given set. We construct a partially ordered set L designed to solve our problem. A typical member X of L is a disjoint collection of countably infinite subsets of A . In detail: X is a set, and each of its members is a countably infinite subset B_i of A . Two different members of X , say B_i and B_j , are required to be disjoint. We partially order L just by set-theoretic inclusion. Any chain in L has an upper bound: Just form the set-theoretic union of all the constituents of the chain. So by Zorn's lemma L has a maximal element, say $Y = \{C_i\}$. If $\cup C_i \neq A$, we

are done. If $\cup C_i$ falls short of A , its complement in A must be finite, for otherwise we could enlarge Y . But this finite ragged remainder can be harmlessly absorbed in one of the C_i 's.

THEOREM 13. *Any infinite set A can be written as a disjoint union $A = B \cup C$, where B , C , and A all have the same cardinal number.*

Proof: By Theorem 12, we express A as a disjoint union of countably infinite subsets S_i . Each S_i can be written as a disjoint union $S_i = T_i \cup V_i$, where T_i and V_i are countably infinite (for instance, we can use the decomposition of the positive integers into even and odd). It suffices to take $B = \cup T_i$, $C = \cup V_i$.

THEOREM 14. *Let d and e be cardinal numbers such that $d \leq e$ and e is infinite. Then $d + e = e$.*

Proof: Obviously $e \leq d + e$. By Theorem 9 it will suffice to prove that $d + e \leq e$. Now $d + e \leq e + e$, and by Theorem 13, $e + e = e$.

EXERCISES

- (Subtraction of cardinals). Let d and e be given cardinal numbers. Prove that the solutions of $d + d_1 = e$ are given as follows:
 - for $d > e$ there is no solution,
 - for $d = e$ and d infinite, d_1 can be any cardinal number $\leq d$,
 - for $d < e$ and e infinite, $d_1 = e$ is the unique solution.
- Prove that $c + c = c$, where c is the cardinal number of the continuum, without using Theorem 13 or Theorem 14. (*Hint:* See Exercise 1 in Section 2.2.)

2.5 CARDINAL MULTIPLICATION

For any cardinal numbers d and e we define their product de to be the cardinal number of the Cartesian product $D \times E$, where $\alpha(D) = d$, $\alpha(E) = e$.

The product of infinite cardinals turns out to be just as trivial as their sum (Theorem 16). Theorem 15 is of course a special case of Theorem 16, but it is convenient to prove it first.

THEOREM 15. *$dd = d$ for any infinite cardinal d .*

Proof: We begin by noting that $\aleph_0 \aleph_0 = \aleph_0$. This follows from Theorem 1 (see also Exercise 2 in Section 2.1). What we are doing in Theorem 15 can be regarded as a generalization of Theorem 1 to arbitrary infinite cardinal numbers.

Let D be a set with $o(D) = d$. Our task is to construct a one-to-one correspondence between D and $D \times D$. We shall do this by Zorn's lemma, and therefore proceed to build the appropriate partially ordered set consisting of all "intermediate stages" of the construction.

Let L be the set of all pairs (E_i, f_i) where E_i is a subset of D and f_i is a one-to-one mapping of E_i onto $E_i \times E_i$. Notice that such an E_i is infinite (if it has more than one element), and that we do have infinite subsets E_i of the required kind, since $\aleph_0 \aleph_0 = \aleph_0$. We declare that $(E_i, f_i) \leq (E_j, f_j)$ if $E_i \subset E_j$ and f_i is the restriction of f_j to E_i . This makes L a partially ordered set in which every chain has an upper bound. Zorn's lemma therefore assures us of the existence of a maximal pair (E, f) .

Before we continue the proof, it is important that we note that E need not be all of D . In somewhat the same fashion as in the proof of Theorem 12, there is a possibility of having a "ragged remainder" which needs special attention; indeed this point is a major technical difficulty of the proof. To illustrate the point, suppose that E is virtually all of D , the complement of E in D consisting of just one element. There is no possibility of extending f to a one-to-one mapping of D onto $D \times D$, for we would have to pair the one remaining element of D with an infinite number of elements (the border of the square in Figure 9). In the same way, if E is uncountable and $D - E$ is countable, f cannot be extended. What we shall have to do is argue that either E is so close to D that our job is finished, or else there is room in the complement for f to be extended.

Write E' for the complement of E in D , and let e and e' be the cardinal numbers of E and E' .

We claim that $e' < e$. For otherwise by Theorem 11 we have $e' \geq e$. This means that we can find within E' a set F with $o(F) = e$. We now enlarge $E \times E$ to the "square" $(E \cup F) \times (E \cup F)$ (see Figure 9). Since $f: E \rightarrow E \times E$ is a one-to-one correspondence, we know that $ee = e$. Hence the three new pieces $E \times F$, $F \times E$, and $F \times F$ all have cardinal number e . By Theorem 14 (used twice), the cardinal number of

$$G = (E \times F) \cup (F \times E) \cup (F \times F)$$

is again e . Hence there exists a one-to-one correspondence between F and G . Combining this with the one-to-one correspondence f between E and $E \times E$, we construct a one-to-one map of $E \cup F$ onto $(E \times E) \cup G = (E \cup F) \times (E \cup F)$ which extends f . This contradicts the maximality of the pair (E, f) . We have sustained the claim that e' must be less than e .

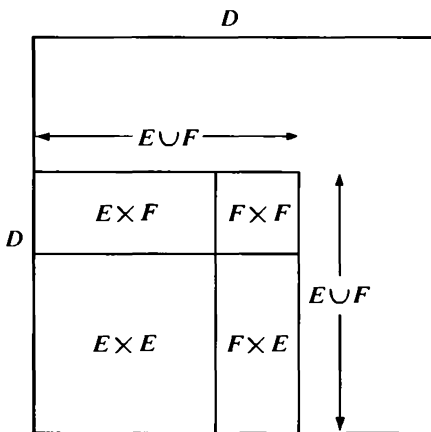


Figure 9

We note next that $d = e + e'$. By Theorem 14, $d = e$. We know that $ee = e$. Hence $dd = d$ has been proved.

THEOREM 16. *Let d and e be cardinal numbers with $d \leq e$, $d \neq 0$, and e infinite. Then $de = e$.*

Proof: Since $d \neq 0$ we have $e \leq de$. On the other hand $de \leq ee$, and $ee = e$ by Theorem 15. We quote Theorem 9.

Consider the special case of the cardinal number c . The fact that $cc = c$ can be interpreted as follows: The number of points in the Euclidean plane is the same as the number of points on a line. When this was discovered by Cantor it was considered very surprising. Possibly, after a century, the shock has lessened. At any rate, the reason it seems paradoxical to us is that we cannot resist attributing to the proposed mapping between the line and the plane some attribute of continuity. It is a fact (going a little beyond the topology treated in this book) that the mapping cannot be continuous in either direction.

Peano did invent a continuous map going from the line onto the entire plane. To be sure, it is very far from one-to-one.

EXERCISES

1. Let $\{A_i\}$, $i = 1, 2, \dots$, be a countable infinity of infinite sets all having the same cardinal number d . Prove that $\cup A_i$ has cardinal number d .

2. (Division of cardinals.) For given cardinal numbers d and e , study the solutions of $dd_1 = e$. (Compare Exercise 1 in Section 2.4.)

2.6 CARDINAL EXPONENTIATION

Given two cardinal numbers d and e , we now propose to attach a meaning to d^e . As motivation, let us examine the finite case. We are initially tempted to think of 3^5 as a 3 multiplied by itself 5 times. We could generalize to the infinite case in this style, but there is another point of view that works better. Let A be a set with 5 elements and B a set with 3 elements. How many functions are there from A to B ? For each element of A we have 3 choices for its image. This selection is made independently 5 times, so we get the answer $3 \cdot 3 \cdot 3 \cdot 3 \cdot 3 = 3^5$. We are ready for the general definition.

DEFINITION. Let d and e be nonzero cardinal numbers. Let D and E be sets with $\alpha(D) = d$, $\alpha(E) = e$. We define d^e to be the cardinal number of the set of all functions from E to D .

Remark: We have assumed d and e to be nonzero so as to avoid any hair-splitting about the number of functions to or from the null set.

The case where either d or e is 1 is uninteresting—see Exercise 1. Also there is nothing new when the exponent is finite—see Exercise 2. Novel thoughts concerning d^e begin with e infinite and $d \geq 2$.

It is desirable to have a symbol for the set of functions from a set E to a set D . The notation D^E is widely used, and it fits well with the notation d^e for exponentiating cardinal numbers. However, we shall use $H(E, D)$, in the same spirit as the basic Hom of homological algebra.

We have used the notation 2^d twice. In Section 2.3 it meant the cardinal number of the power set $P(D)$, where D is a set with cardinal number d . In the present section it means the cardinal number of the set of functions from D to a two-element set. It is easy to identify the two versions, the relevant tool being the characteristic function of a set (compare Exercise 14 in Section 1.4). Let us write T for the set $\{0, 1\}$. Given a subset A of D we attach to it its characteristic function ϕ_A , defined by $\phi_A(x) = 1$ for $x \in A$ and $= 0$ for $x \notin A$. Sending A to ϕ_A is a map of $P(D)$ into $H(D, T)$ which is one-to-one and onto by routine arguments. So $P(D)$ and $H(D, T)$ have the same cardinal number. The cardinal number of $H(D, T)$ is 2^d by definition. In retrospect we see that 2^d was a reasonable notation for the cardinal number of $P(D)$.

At this point let us redeem the promise made in Section 2.3 to prove

that $2^{\aleph_0} = c$, where c is the cardinal number of the set of real numbers. We prove this with the aid of Theorem 9, arguing first that $2^{\aleph_0} \leq c$ and then that $c \leq 2^{\aleph_0}$. Think of the real numbers in their decimal expansions. Fix two digits (other than 0 and 9, so as to avoid the annoying ambiguity that occurs in decimal expansions). Say for definiteness that we pick 5 and 7. The set of real numbers of the form

$$.a_1a_2a_3 \cdots$$

where each a_i is 5 or 7, is in one-to-one correspondence with the set of all functions from the positive integers to the set $\{5, 7\}$. This gives us 2^{\aleph_0} distinct real numbers, proving that $c \geq 2^{\aleph_0}$.

To prove the reverse inequality, we find it convenient to switch to the representation of real numbers in the base 2. (If we took time out at this point to verify that $10^{\aleph_0} = 2^{\aleph_0}$, this change would be unnecessary.) Then every real number gets sent into a sequence

$$b_n \cdots b_2b_1.a_1a_2a_3 \cdots$$

of 0's and 1's. Not all sequences show up. If we agree to exclude decimals that end in a string of 1's, these will not be in the image. For proving that $c \leq 2^{\aleph_0}$, this restriction works in our favor. Likewise the fact that the representation of a real number goes only a finite distance to the left of the decimal point improves our estimate. (This point could be bypassed by working instead with the unit interval; see Exercise 2 in Section 2.2.) Thus $c \leq 2^{\aleph_0}$ is established. Since we have the inequality both ways, we have proved that $c = 2^{\aleph_0}$.

We turn to the consideration of algebraic identities that hold for the operations on cardinal numbers that we have introduced. This might already have been done for addition and multiplication. The identities in question are the commutative and associative laws of addition and multiplication, and the distributive law cutting across them. However, since addition and multiplication of two infinite cardinal numbers give merely the larger of the two, there is very little excitement generated in that discussion. A slightly brighter spark is ignited by the identities involving exponentiation. There are three of them:

$$(10) \quad (d_1d_2)^e = d_1^e d_2^e,$$

$$(11) \quad d^{e_1+e_2} = d^{e_1} d^{e_2},$$

$$(12) \quad (d^e)^f = d^{ef}.$$

It is true that (10) and (11) also trivialize for infinite cardinals. For instance, if d_2 is infinite and $d_1 \leq d_2$, then both sides of (10) are equal to d_2^e . Nevertheless there is perhaps some merit in giving a direct elementary proof of (10), thereby incidentally avoiding the use of Zorn's lemma. Moreover, this can serve as a prelude to the proof of (12), for which no trivialization

is available. The proof of (11) in this elementary style is left for the reader (see Exercise 3).

To prove (10) we select sets D_1, D_2, E with cardinal numbers d_1, d_2, e . With the H notation introduced above, the left side is the cardinal number of $H(E, D_1 \times D_2)$, and the right side is the cardinal number of $H(E, D_1) \times H(E, D_2)$. The discussion will be facilitated by a belated introduction of the projections onto the components of a Cartesian product. Let A and B be any sets. The map $A \times B \rightarrow A$ obtained by sending (a, b) into a is the projection on the first coordinate. We shall write it π_1 , and π_2 will similarly denote the projection on the second coordinate. We can now formally describe the desired mapping from $H(E, D_1 \times D_2)$ to $H(E, D_1) \times H(E, D_2)$. Given $f: E \rightarrow D_1 \times D_2$, we have that the ordered pair $(\pi_1 f, \pi_2 f)$ is an element of $H(E, D_1) \times H(E, D_2)$. That $f \rightarrow (\pi_1 f, \pi_2 f)$ is one-to-one and onto is quite routine.

To prove (12), we let A, B , and C be sets with cardinal numbers d, e , and f , respectively. We have to compare the sets

$$H[C, H(B, A)] \quad \text{and} \quad H(B \times C, A).$$

Let us denote them, for brevity, by LS and RS . A member f of RS is a function from $B \times C$ to A ; thus for every $b \in B$ and $c \in C$ it determines an element $f(b, c)$ of A . From f we have to be able to extract an element g of LS ; g has to be meaningful on C and assign to $c \in C$ a member $g(c)$ of $H(B, A)$. In other words, for $b \in B$ we have to name an element of A to be the value $g(c)(b)$ of $g(c)$ at b . The obvious choice for $g(c)(b)$ is $f(b, c)$. With this choice a mapping $f \rightarrow g$ of RS to LS has been defined. That it is one-to-one is quite routine. To see that it is onto, we have to take a member h of LS and find a member f of RS from which it comes. This is done by taking $f(b, c)$ to be $h(c)(b)$.

Here is an application of (12). Let n be a finite cardinal number, $n \geq 2$. Then $n^{\aleph_0} = c$. For $n^{\aleph_0} \geq 2^{\aleph_0} \geq c$. On the other hand, pick an integer r with $2^r \geq n$. Then $n^{\aleph_0} \leq (2^r)^{\aleph_0} = 2^{r \aleph_0} = 2^{\aleph_0} = c$. (This is one of the entries in the table in Exercise 5.)

As an application of (11), we revisit the equation $cc = c$, which is a special case of Theorem 15 or Theorem 16. We can now handle it as follows:

$$cc = 2^{\aleph_0} 2^{\aleph_0} = 2^{\aleph_0 + \aleph_0} = 2^{\aleph_0} = c.$$

To conclude this section we summarize what can be said about the value of d^e , where, as above, we assume e infinite and $d \geq 2$ to avoid trivialities. The first remark is that for $d \leq 2^e$ we have $d^e = 2^e$ (Exercise 4). Thus attention may be confined to the case $d > 2^e$. Since d^e is in any event greater than e (it is at least 2^e), the interesting thing is to compare d^e with d . There is a "tendency" for d^e to be equal to d ; see for instance Exercise 6. If we look ahead a bit, we can find further information. In the presence

of the continuum hypothesis, Exercise 3 in Section 3.2 shows that $\aleph_2^{\aleph_0} = \aleph_2$. On the other hand, Exercise 2 in Section 3.4 asserts that $\aleph_\omega^{\aleph_0} > \aleph_\omega$. If the continuum hypothesis is assumed (or if we make the weaker assumption that $2^{\aleph_0} < \aleph_\omega$) we get a case of $d^e > d$ with $d > 2^e$. The full elucidation of the facts is considerably beyond the scope of this book.

EXERCISES

- For any cardinal number d , prove that $1^d = 1$ and $d^1 = d$.
- If d is an infinite cardinal number, and n a finite one, prove that $d^n = d$. (*Hint*: Identify d^n with the product of d with itself n times, and use Theorem 15.)
- * Prove (11).
- If e is an infinite cardinal number and d is a cardinal number satisfying $2 \leq d \leq 2^e$, prove that $d^e = 2^e$.
- Verify the correctness of the following table of cardinal exponentiation. Here m and n are finite cardinals, $n \geq 2$. For instance, the second entry in the third row asserts that $c^{\aleph_0} = c$.

	m	\aleph_0	c
n	n^m	c	2^c
\aleph_0	\aleph_0	c	2^c
c	c	c	2^c

- Let d_1, e be infinite cardinals. Write $d = 2^{d_1}$ and suppose that $d \geq 2^e$. Prove that $d^e = d$.
- If $f: A \rightarrow B$ is a given function, the *graph* of f is the subset of $A \times B$ consisting of all $(a, f(a))$.
 - Prove that the map $a \rightarrow (a, f(a))$ is a one-to-one map of A onto the graph of f .
 - * Let D and E be sets with cardinal numbers d and e . Suppose that d is infinite and that $e \leq d$. Prove that d^e is the number of one-to-one maps of E into D . (*Hint*: Of course d^e is an upper bound. To get the reverse inequality, use part (a).)
- Let D be an infinite set with cardinal number d . Prove that D has 2^d subsets of cardinal number d . (*Hint*: Write D as a disjoint union of two subsets D_1, D_2 of cardinal d . Consider all the subsets of D containing D_1 .)
- * Let D be an infinite set with cardinal number d . Prove that the number of one-to-one maps of D onto itself is 2^d . (*Hint*: Obviously the number is

bounded by d^d , which equals 2^d (Exercise 4). In order to exhibit 2^d maps, divide D into two equal halves, observe that one-to-one maps of one half *into* the other can be augmented to one-to-one maps of D onto itself, and use Exercise 7(b).)

- 10.* Let D be an infinite set with cardinal number d . Let e be a cardinal number $\leq d$. Prove that d^e is the number of subsets of D with cardinal e and also the number with cardinal $\leq e$. (That d^e is an upper bound is easy. To exhibit d^e subsets of cardinal e , note that there are d^e one-to-one maps of a set with cardinal e into D (Exercise 7(b).) By Exercise 9, these fall into equivalence classes of size 2^e , with the number of equivalence classes equal to the number of subsets of cardinal e . If $d^e > 2^e$, argue that the number of equivalence classes must be d^e . If $d^e = 2^e$, then we only need to find 2^e subsets of cardinal e . By Exercise 9 we can even do this inside a set of cardinal e .)
- 11.* Let E and D be sets with cardinal numbers d and e . Assume that e is infinite and that $e \geq d$. Prove that the number of functions from E onto D is d^e .

3

Well-ordering; The Axiom of Choice

3.1 WELL-ORDERED SETS

We turn in this chapter to another basic contribution made by Cantor: his systematic treatment of the process of going on forever again and again, which we discussed informally (and inadequately) in Section 2.3. The pertinent concept is embodied in the following definition.

DEFINITION. A well-ordered set is a chain in which every nonempty subset has a smallest element.

We could equally well define a well-ordered set to be a partially ordered set in which every nonempty subset has a smallest element, for by enforcing the condition on a two-element subset, we make the elements comparable.

What does a well-ordered set C look like? There is a smallest element a_0 in C . In the complement of $\{a_0\}$ there is a smallest element a_1 . This process continues, and produces an ascending sequence

$$a_0 < a_1 < a_2 < \dots$$

which goes on forever, except in the trivial case where C is finite. The a_i 's might exhaust C . Failing this, there is a smallest element in the comple-

ment of $\{a_i\}$ within C . This gets us started in going on forever a second time. And so on. If this seems too mysterious, we can exhibit with very concrete rational numbers a well-ordered set consisting of two such sequences, one "on top" of the other:

$$0, \frac{1}{2}, \frac{3}{4}, \frac{7}{8}, \dots, 1, 1\frac{1}{2}, 1\frac{3}{4}, 1\frac{7}{8}, \dots$$

We shall not at present attempt to peer farther down this misty road. Let us start the theory.

THEOREM 17. *A chain is well-ordered if and only if it does not contain an infinite descending sequence.*

Remark: In modern algebra considerable use is made of the following concept: A partially ordered set satisfies the *descending chain condition* if it contains no infinite descending sequence (or equivalently, every subset has a minimal element). In this language, a well-ordered set is a chain satisfying the descending chain condition.

Proof: If a chain C is well-ordered, it cannot contain an infinite descending sequence, for this would be a subset without a smallest element. Conversely, suppose that the chain C is not well-ordered. Then there is a subset S of C with no smallest element. We pick any a_1 in S . Since a_1 is not minimal in S we can find $a_2 \in S$ with $a_2 < a_1$. Similarly we must have $a_3 \in S$ with $a_3 < a_2$. The sequence $a_1 > a_2 > a_3 > \dots$ is infinite descending.

We introduce two related concepts, giving them the names "ideal" and "segment," in the hope that the exposition will thereby become a little smoother. They might as well be defined in any partially ordered set.

DEFINITION. An *ideal* in a partially ordered set L is a subset I such that $x \in I, y \leq x$ imply $y \in I$, i.e., along with any element, I is to contain all smaller ones.

DEFINITION. Let a be an element in a partially ordered set L . The *segment* $S(a)$ defined by a is the set of all x in L with $x < a$. A subset of L is a segment if it is the segment determined by some element.

Segments and ideals do not differ very much in a well-ordered set.

THEOREM 18. *An ideal in a well-ordered set C is either all of C or a segment of C .*

Proof: Suppose an ideal I is not all of C . Let a be the smallest element of the complement of I . We claim that $I = S(a)$. If $x < a$ then necessarily $x \in I$, since a is minimal in the complement of I . Conversely, if $x \in I$ then $x < a$; for otherwise $x \geq a$ and $a \in I$ since I is an ideal.

In the light of Theorem 18, the purpose of our duplicate terminology is to avoid repeated use of the phrase "a segment of C or all of C ."

Well-ordering need only be checked on segments.

THEOREM 19. *If every segment in a chain C is well-ordered, then C is well-ordered.*

Proof: If C is not well-ordered, then by Theorem 17 it contains a descending sequence $a_1 > a_2 > a_3 > \dots$. Then the segment $S(a_1)$ is clearly not well-ordered.

Recall that in Section 1.4 we defined an order-isomorphism between partially ordered sets to be a mapping f that is one-to-one, onto, and preserves order, i.e. $f(x) \geq f(y)$ if and only if $x \geq y$.

THEOREM 20. *A well-ordered set cannot be order-isomorphic to one of its segments.*

Proof: Suppose on the contrary that the well-ordered set C is order-isomorphic to the segment $S(a)$, $a \in C$, and let $f: C \rightarrow S(a)$ be the alleged order-isomorphism. We of course have $f(a) < a$. Thus the set T of elements x in C with $f(x) < x$ is nonempty. Let b be the smallest member of T . Then $f(b) < b$, and, since f is order-preserving $f(f(b)) < f(b)$. This shows that $f(b)$ is a smaller element of T than b , a contradiction.

The next theorem is the analogue, in the present context, of Theorem 10. Well-ordering is tailor-made to ensure a simple proof.

THEOREM 21. *Let A and B be well-ordered sets. Then either A is order-isomorphic to an ideal in B or B is order-isomorphic to an ideal in A .*

Proof: Call an element a in A "good" if there exists an order-isomorphism of $S(a)$ onto a segment $S(b)$ for some $b \in B$. Note that it follows from Theorem 20 that the element b is unique. Let I be the set of good elements. It is routine to check the following statements: I is an ideal in A , the mapping $f: I \rightarrow B$ defined by $f(a) = b$ (where b is related to a as above) is one-to-one and order-preserving, and the image J of f is

an ideal in B . If both statements (in the conclusion of Theorem 21) fail, let x and y be the smallest elements in the complements of I and J . It is immediate that x is good, since in fact $S(x) = I$, $S(y) = J$.

The following theorem was of immense importance in the historical development of the subject. We insert it quietly at this point, with just the remark that of course we use Zorn's lemma in the proof.

THEOREM 22. *Any set can be well-ordered.*

Proof: The partially ordered set that gets Zornified is simply the set (say L) of all nonempty subsets of the given set A , each equipped with a well-ordering. If B and C are two such objects, by $B \leq C$ we mean that B is an ideal in C (and of course the ordering on B agrees with the ordering of C when restricted to B). We illustrate the idea. All one-element subsets of A appear in L , and there is no choice for the ordering. A two-element subset, say $\{a, b\}$, appears twice: once with the order $a < b$, and once with the other order $b < a$. Call these B_1 and B_2 respectively. Let C be the subset $\{a, b, c\}$ with the order $a < b < c$. Then $B_1 \leq C$, but $B_2 \not\leq C$ does not hold, since the orders do not agree. A finite subset, with n elements, will appear in L precisely $n!$ times. A countably infinite subset will appear many times indeed (compare Exercise 10). Note that we do not know in advance which subsets of A will show up at all (though of course every countable subset surely does). In fact, the whole point of the proof is to establish finally that the whole set A does make an appearance in L .

We claim that every chain in L possesses an upper bound. Let $\{B_i\}$ be a chain in L . The proposed upper bound, as a subset of A , is simply $B = \cup B_i$. We must define an ordering on B , check that it is a well-ordering, and then prove that every B_i is an ideal in B . Let us treat the various points on the agenda under numbered headings.

(1) *Defining the order on B .* Let $x, y \in B$ be given. These two elements got into B on the authority of being in one of the B_i 's. Say $x \in B_i$ and $y \in B_j$. Now B_i and B_j are comparable one way or the other; say $B_i \leq B_j$ for definiteness. Then both x and y lie in B_j . There (if x and y are different) exactly one of the relations $x < y$, $y < x$ holds. What if another investigator decides on the ordering of x and y by looking at a different set, say B_k , containing x and y ? All is well, since B_j and B_k are comparable, and therefore the verdict on the ordering of x and y does not depend on which of the B_i 's is used. We have a well-defined relation on B .

(2) *B is a partially ordered set.* Reflexivity is trivial. So is antisymmetry, for (as is implicit above) if both $x \leq y$ and $y \leq x$ hold, this would

be visible in a single B_k . For transitivity, we hunt up a B_i containing all three of the given elements.

(3) *B is a chain.* This is clear; x and y were found to be comparable in a B_i containing both.

(4) *B is a well-ordered chain.* By Theorem 18 it suffices to show that every segment in B is well-ordered. Take the segment $S(x)$, and suppose that $x \in B_i$. We claim that the entire segment $S(x)$ lies in B_i ; this will do the trick, since we know that B_i is well-ordered. So let y be an element of B with $y < x$. We need to show that $y \in B_i$. We have $y \in B_j$ for some B_j . If $B_j \leq B_i$, all is well. If $B_i \leq B_j$, we observe that B_i is an ideal in B_j , and deduce from $x \in B_i$, $y < x$ that $y \in B_i$.

(5) *B_i is an ideal in B .* We just proved this.

We have sustained the claim that every chain in B has an upper bound. This entitles us to reap the rewards of Zorn's lemma: There is a maximal element in L . Call it D . Recall that D is a subset of A , equipped with a well-ordering. We claim that D must be all of A . For suppose that u is missing in D . We concoct a new well-ordered set out of $\{D, u\}$ by placing u on top of D ($u \geq x$ for all $x \in D$). That $\{D, u\}$ is well-ordered is immediate (see Exercise 6 for a bit of a generalization). Moreover, D is an ideal in $\{D, u\}$. This contradiction of maximality completes the proof of Theorem 22.

It is time to take a retrospective look at Theorem 10. We can now prove it as follows: Well-order A and B by Theorem 22 and then apply Theorem 21. It is a tenable position that this is a more agreeable proof of Theorem 10 than the one given in Section 2.3.

Why, then, was Theorem 10 proved à la Zorn in Section 2.3? Answer: for pedagogical reasons. It was time to discuss the comparability of sets, timely to introduce Zorn, and possible to delay well-ordering. The current state of affairs in mathematics calls for the well-equipped mathematician to be able to use Zorn's lemma, apply well-ordering, or do both or neither, as the occasion requires. Furthermore, other things being equal, everybody (or at any rate nearly everybody) would rather Zornify. And it is a fact that 99% of the time, a proof by Zorn's lemma is slicker than by well-ordering. (However, in Appendix 2 we present an example (one of the 1%) where well-ordering is better.)

We return to the theory of cardinal numbers. The results on well-ordering enable us to enrich our picture of what cardinal numbers look like.

THEOREM 23. *Any set of cardinal numbers is well-ordered (in the ordering we have introduced on cardinal numbers).*

Remarks: 1. It would be playing with fire to talk about “the set of all cardinal numbers,” so we cautiously use the neutral phrase “any set of cardinal numbers.”

2. With Theorem 23 at hand we picture the initial string of infinite cardinal numbers as

$$\aleph_0, \aleph_1, \aleph_2, \dots, \aleph_n, \dots$$

Where does $2^{\aleph_0} = c$ fit in? We shall return to this in Section 3.4.

Proof of Theorem 23: Assume the contrary. Then by Theorem 17 there must exist a strictly descending sequence $d_1 > d_2 > d_3 > \dots$ of cardinal numbers. For each d_i pick a set D_i with $o(D_i) = d_i$. By Theorem 22 we may assume that each D_i is well-ordered. Apply Theorem 21 to D_i and D_1 . In view of the inequality $d_1 > d_i$, it must be the case that D_i is order-isomorphic to an ideal in D_1 that is not all of D_1 , say $S(b_i)$. Moreover, the inequality $d_i > d_{i+1}$ implies that $b_i > b_{i+1}$ must hold. The b_i 's are a forbidden descending sequence in D_1 .

EXERCISES

1. Prove that any subset of a well-ordered set is well-ordered (in the inherited ordering).
2. Prove that in a well-ordered set every subset bounded from above has a least upper bound.
3. Let C be a chain in which every countable subset is well-ordered. Prove that C is well-ordered.
4. (a) Let C be a well-ordered set and let $f: C \rightarrow C$ be a one-to-one function such that $a \leq b$ implies $f(a) \leq f(b)$. Prove that $a \leq f(a)$ for all a in C .
(b) Deduce Theorem 20 from part (a).
5. Let C be a set of real numbers which is well-ordered in its usual ordering. Prove that C is countable. (*Hint:* Insert a rational number between every element of C and the element directly above it.)
6. Let A be a chain. Let B and C be subsets of A with $A = B \cup C$. Suppose that B and C are well-ordered (in the ordering they inherit from A). Prove that A is well-ordered.
7. Extend Exercise 6 to a finite union of well-ordered sets. Is it true for infinite unions?
- 8.* Let X be a countable set. Let C be a collection of subsets of X with the following property: For any ascending sequence $A_1 \subset A_2 \subset A_3 \subset \dots$ of members of C , $\cup A_i$ is again in C . Prove that C has a maximal element.

- 9.* Prove that any chain has a cofinal well-ordered subset. (For the definition of *cofinal*, see Exercise 5 in Section 2.1.)
- 10.* Let D be an infinite set with cardinal number d . Prove that the following all have cardinal number 2^d :
- The number of ways of partially ordering D .
 - The number of ways of making D a chain.
 - The number of ways of well-ordering D .
- 11.* Let C be a chain which every subset has a top element or a bottom element (or both). Prove that C consists of a well-ordered set surmounted by the dual of a well-ordered set. (The dual of a partially ordered set is the partially ordered set obtained by reversing the ordering.)

3.2 ORDINAL NUMBERS

We proceed to effect a passage from well-ordered sets to ordinal numbers. This transition is entirely analogous to the passage from general sets to cardinal numbers.

We attach to every well-ordered set an *ordinal number*; two well-ordered sets are awarded the same ordinal number if and only if they are order-isomorphic. We typically use the letters α, β, \dots for ordinal numbers.

We introduce an ordering on ordinal numbers. Let α, β be ordinal numbers and let A, B be well-ordered sets representing them. We write $\alpha \leq \beta$ if A is order-isomorphic to an ideal in B . The definition is easily seen to be independent of the choice of A and B . Reflexivity and transitivity are immediate. Antisymmetry follows from Theorem 20; so we have a partial ordering. By Theorem 21 it is a linear ordering.

THEOREM 24. *Let α be any ordinal. The ordinals β with $\beta < \alpha$ form a well-ordered set with ordinal number α .*

Proof: Fix a well-ordered set A with ordinal number α . Given $\beta < \alpha$, take a well-ordered set B representing β . Then B is order-isomorphic to a segment in A , say $S(b)$. It is straightforward to check that the mapping $\beta \rightarrow b$ is an order-isomorphism onto A .

In view of Theorem 19 we have a corollary.

THEOREM 25. *Any set of ordinal numbers is well-ordered.*

As in Theorem 23, we avoid the phrase "the set of all ordinal numbers." This time the matter is worth exploring a little further, for the act of

taking all ordinal numbers resembles respectable mathematics more than taking the set of all nonself-members (as in Russell's paradox). We present the Burali-Forti paradox.

Cesare Burali-Forti (1861–1931) taught at the Royal Academy of Artillery and Engineering in Turin from 1903 to 1931. His famous paper appeared in the 1897 *Atti Torino*.

Let W be the set of all ordinal numbers. By Theorem 25, W is well-ordered. Its ordinal number, say γ , is thus a member of W . By Theorem 24, W is order-isomorphic to $S(\gamma)$. This contradicts Theorem 20.

The moral is clear: We must treat large sets of ordinal numbers with delicacy. How to do this is part of the subject called *axiomatic set theory*, and lies beyond the scope of this book.

No further theorems will be proved in this section. We shall collect a number of observations which we hope will help the reader cope with ordinal numbers when he encounters them in the literature.

(1) *Ordinals to cardinals and back.* If two well-ordered sets are order-isomorphic they obviously have the same cardinal number. Hence we have a mapping from ordinal numbers to cardinal numbers. This mapping is onto by Theorem 22. It is not one-to-one (see observation 2). However, there is a unique *smallest* ordinal number with cardinal d . It is called the *initial ordinal* of cardinal d .

(2) *Notation.* The ordinals start with 0 (the ordinal number of the empty set). Then we move through the finite ordinals

$$(13) \quad 0, 1, 2, 3, \dots, n, \dots$$

We are using the same notation as is used for the finite cardinals, since there is no point in making a distinction. The well-ordered set (13) has for its ordinal number the first infinite ordinal. Cantor gave it the notation ω (omega). It is the initial ordinal of the cardinal \aleph_0 . The next well-ordered set is

$$(14) \quad 0, 1, 2, \dots, n, \dots, \omega.$$

The notation for the ordinal number of (14) is $\omega + 1$. (The notation suggests the possibility of our defining addition on ordinal numbers as well as cardinal numbers. We shall not pursue this, except for a little exploration in exercises.) Note that $\omega + 1$ is also a countable ordinal (by a *countable ordinal* we mean one corresponding to the cardinal \aleph_0). The ordinals ω and $\omega + 1$ are different; for instance (14) has a last element while (13) does not.

We call an ordinal number (different from 0) a *limit ordinal* if it has no immediate predecessor. ω is the first limit ordinal. The second is 2ω , the ordinal number of

$$0, 1, \dots, n, \dots, \omega, \omega + 1, \dots, \omega + n, \dots$$

There follow the limit ordinals $3\omega, 4\omega, \dots$. The ω th limit ordinal is ω^2 . After that we have $\omega^3, \omega^4, \dots$ and the end of this road is an ordinal written ω^ω . (Since I once spent a sleepless night over this, let me mention a possible pitfall. The cardinal of ω^ω is not $\aleph_0^{\aleph_0} = c$. It is still a countable ordinal.) There is a natural meaning for

$$\omega^{\omega^\omega}, \omega^{\omega^{\omega^\omega}}, \dots$$

and at the end of this trip we arrive at an ordinal which Cantor called ϵ . This is still a countable ordinal! In fact we have just skimmed the surface of the countable ordinals.

Remark: As the symbols 2ω and ω^ω suggest, there are useful definitions of multiplication and exponentiation of ordinal numbers. We shall not touch on this at all.

Since cardinals are well-ordered, it is natural to write them as \aleph 's with an ordinal subscript. It is customary to begin this with the infinite cardinals:

$$\aleph_0, \aleph_1, \aleph_2, \dots, \aleph_n, \dots, \aleph_\omega, \aleph_{\omega+1}, \dots$$

We have the notion of *limit cardinal*: a cardinal with no immediate predecessor. \aleph_0 is the first limit cardinal, \aleph_ω the next, $\aleph_{2\omega}$ the third, \dots

The initial ordinal for the cardinal \aleph_1 is written Ω (capital omega). Thus Ω is the first uncountable ordinal. Preceding it we have the finite ordinals and then \aleph_1 countable ordinals. The ordinals up to Ω , or up to and including Ω , form a popular starting point for bizarre examples in topology.

(3) *Bouncing back and forth.* Start with the ordinal 0. Take the \aleph with that subscript: \aleph_0 . Pass to the corresponding initial ordinal: ω . Then we have the cardinal \aleph_ω , then the initial ordinal of \aleph_ω , etc. In very short order we reach some very large looking cardinals and ordinals.

(4) *Transfinite induction.* Let $P(\lambda)$ be a statement about the ordinal λ . Suppose $P(0)$ is true, and suppose that we know $P(\alpha)$ to be true if $P(\beta)$ holds for all $\beta < \alpha$. Then $P(\lambda)$ is true for all λ . For if there is an ordinal for which P is false, then there is a smallest one, a contradiction.

This method of proof is called *transfinite induction*. It is quite analogous to an induction over the positive integers. There is one difference: Most of the time it turns out to be necessary to distinguish two cases in the inductive step, according as α is, or is not, a limit ordinal.

(5) *Transfinite construction.* When ordinal numbers appear in the literature the context is very often one of constructing certain objects, one for each ordinal. The difference between this and transfinite induction is really just psychological, since we could rephrase the construction as an existence theorem.

See Appendix 2 for an application of transfinite construction in the theory of abelian p -groups.

EXERCISES

1. Prove that $\aleph_1^{\aleph_0} = 2^{\aleph_0}$. (Observe that $\aleph_1 \leq 2^{\aleph_0}$ and use Exercise 5 in Section 2.6.)
- 2.* The sum $\alpha + \beta$ of ordinal numbers α, β is defined as follows: Pick disjoint well-ordered sets A and B representing α and β , and order $A \cup B$ by making every element of B exceed every element of A . After proving that $A \cup B$ is well-ordered (compare Exercise 6 in Section 3.1), let $\alpha + \beta$ be its ordinal number.
 - (a) Is $\alpha + \beta$ always equal to $\beta + \alpha$?
 - (b) Prove that $\alpha + \beta = \alpha + \gamma$ implies $\beta = \gamma$.
 - (c) Give an example where $\beta + \alpha = \gamma + \alpha$ fails to imply $\beta = \gamma$.
- 3.* Prove that $\aleph_2^{\aleph_0} = \max(\aleph_2, 2^{\aleph_0})$. (Well-order a set A of cardinal \aleph_2 with the initial ordinal of that cardinal. Any map of a countable set into A has range bounded by an ordinal of cardinal $\leq \aleph_1$. There are \aleph_2 of these ordinals. Hence we get the estimate

$$\aleph_2^{\aleph_0} \leq \aleph_2 \aleph_1^{\aleph_0}.$$

Use Exercise 1.)

- 4.* Let L be a partially ordered set satisfying the descending chain condition (that is, L contains no infinite descending sequence). Prove that the order in L can be strengthened so as to make L a well-ordered chain (compare Exercise 5 in Section 2.3).

3.3 THE AXIOM OF CHOICE

Up to this point we have taken the attitude which is standard in the current mathematical literature: If a theorem comes up which seems to need a "transfinite method," go ahead and prove it, preferably using Zorn's lemma. In this section we turn back to the foundations of set theory, and ask whether Zorn's lemma might itself be proved from something more plausible.

As we remarked in Section 2.3, Zermelo proposed an axiom for this purpose.

THE AXIOM OF CHOICE. *Let A be any set. Then there exists a function f , defined on the nonempty subsets of A and taking values in A , which assigns to every such subset one of its members.*

Ernst Zermelo's important paper appeared in *Math. Annalen* **59** (1904), 514–6. It was actually an extract from a letter to Hilbert, and in the final paragraph Zermelo thanks Erhard Schmidt for suggesting the idea. In the very next volume of the *Annalen* (pp. 194–5), it is interesting to read Borel's gloomy evaluation of the axiom of choice. Zermelo returned with a spirited defense in the *Annalen* **65** (1908), 107–128.

Let us line up three statements:

1. The axiom of choice (AC).
2. Zorn's lemma (ZL).
3. Any set can be well-ordered (WO).

Now before saying anything further, let it be noted that the three statements are equivalent (as we shall shortly show). Nevertheless, most mathematicians take different attitudes to the three statements. Let me record my (purely subjective) feelings. AC is utterly acceptable: Just pick an element in each set! ZL is somewhat less acceptable. I would hate to accept WO as an axiom.

Mathematicians' feelings about AC are reflected in their work. Some steer clear of AC and work in "safe" corners of mathematics. At the other extreme are those who make no distinction between mathematics with AC and mathematics without it. (My middle position: I try to remember to make a note when I use it, but I do not hesitate to use it.)

We shall prove the equivalence of the three statements by going in the circle $AC \Rightarrow ZL \Rightarrow WO \Rightarrow AC$. But before doing so, we remark that it would fit the history of the subject better if the plan were $AC \Rightarrow WO \Rightarrow ZL \Rightarrow AC$. Done in this style, $WO \Rightarrow ZL$ would be a transfinite induction (see the preceding section for hints on how to carry out this enterprise), $ZL \Rightarrow AC$ would be a fairly straightforward Zornification, and $AC \Rightarrow WO$ would be Zermelo's proof, very similar to the second proof we shall give that $AC \Rightarrow ZL$.

Two-thirds of the project $AC \Rightarrow ZL \Rightarrow WO \Rightarrow AC$ can be disposed of very quickly: $ZL \Rightarrow WO$ is Theorem 22, and for $WO \Rightarrow AC$ all we have to do is let the choice function take, as its value on a set, the smallest element of that set. However, the remaining third ($AC \Rightarrow ZL$) will occupy our attention for quite a while.

We make a small change in strategy. Let us recall the full statement of ZL.

ZL: Let L be a partially ordered set in which every chain has an upper bound. Then L has a maximal element.

Compare this with ZL'.

ZL': Let L be a partially ordered set in which every chain has a least upper bound. Then L has a maximal element.

Obviously $ZL \Rightarrow ZL'$. Now (in the first of the two proofs below) we are actually going to prove that $AC \Rightarrow ZL'$. So after that we shall have to supply $ZL' \Rightarrow ZL$. In fact, let us do that right now. Sad to tell, we do this via a third version of ZL !

ZL'' : Any partially ordered set contains a maximal chain.

(This is the version of Zorn's lemma which appeared in Hausdorff's book, as mentioned in Section 2.3.)

Since $ZL \Rightarrow ZL'$ is obvious (as noted above), we are left with $ZL' \Rightarrow ZL''$ and $ZL'' \Rightarrow ZL$.

Proof of $ZL' \Rightarrow ZL''$: Take the set of all subchains of the given partially ordered set. We order by inclusion, and we get a partially ordered set satisfying the hypothesis of ZL' (the required least upper bound is simply the set-theoretic union). Thus we get a maximal chain.

Proof of $ZL'' \Rightarrow ZL$: On the authority of ZL'' , we select a maximal chain C in L . Let u be an upper bound of C . u must be maximal, for if $v > u$ then adjoining v to C yields a larger chain.

Query: Why did we take our basic version of Zorn's lemma in the form ZL in preference to the form ZL' ? The answer is that it is a matter of convenience, or taste, or both. It is a fact that everywhere we used Zorn's lemma the upper bound that came up was in fact a least upper bound. This is almost always the case in the literature. So there is very little difference in practice between ZL and ZL' . At any rate, by using ZL we are spared the trouble of checking (for about half a second) that the pertinent chains have least upper bounds, and not just upper bounds.

We shall give two proofs of $AC \Rightarrow ZL$. In the first, as remarked above, we actually prove $AC \Rightarrow ZL'$. The first proof might be preferred on expository grounds, since the very notion of well-ordering gets bypassed. The second uses well-ordered sets as a *tool* (note that it is out of the question to use the theorem that any set can be well-ordered), and follows a path that might be considered more natural.

The first proof is preceded by two lemmas. The second of these lemmas is obviously stronger than the first, but it makes for simpler exposition to break the proof up in this way.

LEMMA 1. *Let L be a partially ordered set in which every chain has a least upper bound. Then there cannot exist a function $g: L \rightarrow L$ such that for every x in L , $g(x)$ is directly above x (i.e., $g(x) > x$, and no element of L lies strictly between $g(x)$ and x).*

LEMMA 2. *Let L be a partially ordered set in which every chain has a least upper bound. Then there cannot exist a function $g:L \rightarrow L$ satisfying $g(x) > x$ for all x in L .*

Remarks: 1. Zorn's lemma, in the form ZL or ZL', would tell us that L has a maximal element and give us the conclusion of Lemma 2 at once. But of course we must not use Zorn's lemma; we are engaged in proving it.

2. On the other hand, it would be acceptable (at least for our present purposes) to use AC in proving the lemmas. But it is nice that the proof does not need it.

Proof of Lemma 1: We assume the existence of g and head for a contradiction.

Fix an element $z_0 \in L$. If L_0 denotes the set of all a in L with $a \geq z_0$, then all our hypotheses apply to L_0 . Thus we can work exclusively in L_0 and we do so; we further simplify notation by switching from L_0 to L . So: We have arranged that L has a unique bottom element z_0 .

Following the terminology of Halmos on p. 64 of *Naive Set Theory* we call a subset A of L a *tower* if it has the following three properties:

- (a) $z_0 \in A$,
- (b) $g(A) \subset A$,
- (c) Along with any chain, A contains the least upper bound of the chain.

Towers exist: For instance L is a tower. Let T be the intersection of all towers. Clearly T is again a tower. The plan of the proof is to show that T is a chain. Then the least upper bound t of T lies in T , and $g(t)$ lies in T , contradicting $g(t) > t$. So the proof is finished the moment we show that T is a chain.

Let V be the set of elements in T which are comparable with every element of T . In other words, $b \in V$ means that $b \in T$ and that for any u in T either $b \leq u$ or $u \leq b$. If we show that $V = T$ then T is a chain and all is done. Since $V \subset T$, we need only prove $V \supset T$, and we shall do this by proving that V is a tower. Property (a) is clear. We next verify (c). Assume that $\{b_i\}$ is a chain in V and let u be in T . Let b be the least upper bound of $\{b_i\}$. If a single b_i satisfies $b_i \geq u$, then $b \geq u$. If on the other hand every $b_i \leq u$ then $b \leq u$ (by the defining property of a least upper bound). Thus b is comparable with u and hence $b \in V$.

It remains to check (b), i.e., given $b \in V$ we must prove $g(b) \in V$. For this purpose we form W , the set of all elements w in T satisfying either $w \leq b$ or $g(b) \leq w$. We claim W is a tower and again we have to check (a) (c).

(a) Since $z_0 \leq b$ we have $z_0 \in W$.

(b) Take $w \in W$. We have to prove that $g(w) \in W$. Of course, $g(w) \in T$, so we have only to check that $g(w) \leq b$ or $g(b) \leq g(w)$.

Case I: If $g(b) \leq w$ then $g(b) < g(w)$, and all is well.

Case II: Assume that $w \leq b$. Since $b \in V$ and $g(w) \in T$ we have $g(w) \leq b$ or $b \leq g(w)$. The possibility $g(w) \leq b$ is acceptable, so we assume $b \leq g(w)$. Together we then have $w \leq b \leq g(w)$. By the hypothesis that no element of L lies strictly between w and $g(w)$ we have $b = w$ or $g(w) = b$. Now $g(w) = b$ is even better than $g(w) \leq b$, and $w = b$ implies $g(w) = g(b)$, better than $g(w) \geq g(b)$. We have checked all possibilities and have proved $g(w) \in W$.

(c) Let $\{w_i\}$ be a chain in W and let y be its least upper bound. We have to prove $y \in W$, that is, we have to prove $y \leq b$ or $g(b) \leq y$. We know that each w_i satisfies $w_i \leq b$ or $g(b) \leq w_i$. If every $w_i \leq b$ then $y \leq b$. Otherwise we must have $g(b) \leq w_i$ for at least one w_i , and then $g(b) \leq y$.

We have completed the proof that W is a tower. It follows that $W = T$. This means: For every u in T we have $u \leq b$ or $g(b) \leq u$. Since $u \leq b$ implies $u \leq g(b)$ we get that $g(b)$ is comparable with u . This is the criterion for $g(b)$ to lie in V . With this we have completed the proof that V is a tower and, as noted above, we have proved Lemma 1.

Proof of Lemma 2: We prove Lemma 2 by switching the scene of activity to a suitable partially ordered set of subsets of L . Let \mathfrak{L} be the set of chains which are contained in L . Partially order \mathfrak{L} by inclusion. Define $h: \mathfrak{L} \rightarrow \mathfrak{L}$ as follows: If $C \in \mathfrak{L}$ and has a top element x , set $h(C) = C \cup g(x)$; if $C \in \mathfrak{L}$ and has no top element, set $h(C) = C \cup$ (the least upper bound of C). In either case $h(C)$ is larger than C by exactly one element. \mathfrak{L} is clearly a partially ordered set in which every chain has a least upper bound. Relative to \mathfrak{L} and h (in place of L and g) we fulfill the hypothesis of Lemma 1. We have thereby concluded the proof of Lemma 2.

We are now ready for the long promised $AC \Rightarrow ZL'$.

Proof of $AC \Rightarrow ZL'$: For $x \in L$, let M_x be the set of all y in L with $y > x$. If any M_x is empty we have our maximal element. So assume that every M_x is nonempty. Let f be a choice function on L , as provided by AC. Define $g: L \rightarrow L$ by $g(x) = f(M_x)$. Then $g(x) > x$ for all x . We have contradicted Lemma 2.

We turn to the second proof that was promised.

Second proof that AC implies Zorn's lemma: This time we are taking Zorn's lemma in the form ZL. So let L be a partially ordered set in which every chain has an upper bound. We are to prove that L has a maximal element. Assume the contrary.

We first observe that any chain C in L has an upper bound that is not in C . For let u be any upper bound of C . u is not maximal, so $u < v$ for some v in L , and v is an upper bound for C with $v \notin C$.

Let f be a choice function for the set L . For any chain C , let C^* denote the set of upper bounds of C that are not contained in C ; we have just observed that C^* is not empty. Define $g(C) = f(C^*)$. Thus g is a function defined on every chain in L and taking as value an upper bound which is not an element of the chain.

We fix an element $z_0 \in L$ for the rest of the proof. Let B be a subset of L , equipped with a well-ordering. We shall call B "special" if z_0 is the smallest element of B and furthermore, for any y in B other than z_0 , $y = g[S(y)]$. (Note that $S(y)$ is the segment of B determined by y .)

The crucial thing that needs to be proved is this: If B and D are special well-ordered subsets of L , then either B is an ideal in D or D is an ideal in B . To see this, let E be the set of elements u with the following properties: $u \in B \cap D$ and $S(u)$ is the same, whether computed in B or in D . It is immediate that E is an ideal in both B and D . Suppose first that E is a proper subset of both B and D . Then $E = S(v)$ in B and $E = S(w)$ in D . Since B and D are special well-ordered subsets of L , v and w are both equal to $g(E)$. Now that we know $v = w$, it follows that v lies in E , a contradiction. We conclude that E must coincide with all of B or D , say $E = B$. Thus B is an ideal in D .

From this point on, the proof closely resembles the corresponding portion of the proof of Theorem 22. We form the set-theoretic union, say G , of all special well-ordered subsets of L and argue easily that G admits a natural well-ordering which is again special. But we can construct a larger special well-ordered set by placing $g(G)$ on top of G . This contradiction completes the proof of $AC \Rightarrow ZL$.

We shall briefly mention two alternative versions of AC, labelling them AC' and AC'' .

AC' : Let $\{A_i\}$ be a collection of disjoint nonempty sets. Then there exists a set B having exactly one element in common with each A_i , and containing no other elements.

To see that $AC \Rightarrow AC'$, form $C = \cup A_i$, take a choice function f for C , and let B be the set whose members are $f(A_i)$. Conversely, assume AC' and let D be a set. The difficulty in applying AC' to get a choice function for D is that the nonempty subsets of D overlap. But this is easily remedied

by the device of taking ordered pairs, as we did in defining the sum of two cardinal numbers (Section 2.2).

AC'': The Cartesian product of nonempty sets is nonempty.

We leave it to the reader to fit AC'' into the picture.

In concluding this section we recall that something we called the "countable axiom of choice" was informally mentioned several times. It can be taken to be AC' for a countable collection of sets, or AC'' for a Cartesian product of countably many sets.

EXERCISES

1. Let A and B be sets admitting a function $f: A \rightarrow B$ which is onto. Prove that $\alpha(B) \leq \alpha(A)$. (*Remark:* This could undoubtedly have been slipped in at the beginning of Section 2.3 without attracting undue attention. But it does require AC in a very direct way. Compare with Exercise 1 in Section 2.1.)
2. Let D be an infinite set with cardinal number d . Prove that d is the number of finite subsets of D . (*Hint:* Observe that there is a natural map of the Cartesian product of D with itself n times onto the subsets of D with n or fewer elements. Use the preceding exercise and then Exercise 1 in Section 2.5.)
3. Prove that any chain in a partially ordered set can be expanded to a maximal chain.
4. Prove that any partially ordered set which is not a chain contains at least two maximal chains.
5. Prove the following directly from the axiom of choice: For any function f from a set X to itself there exists $g: X \rightarrow X$ with $fgf = f$.
6. Prove the converse of Exercise 5. In detail: Suppose that for any set X and any function $f: X \rightarrow X$ there exists $g: X \rightarrow X$ with $fgf = f$. Prove that any set has a choice function. (*Hint:* Use the version AC'.)

3.4 THE CONTINUUM PROBLEM

In the well-ordered array of infinite cardinals

$$\aleph_0, \aleph_1, \aleph_2, \dots$$

where does $c = 2^{\aleph_0}$ fit in? It is at least \aleph_1 , of course. Cantor conjectured that $c = \aleph_1$ and this is called the *continuum hypothesis*. Hilbert made this the first of his 23 famous problems.

There was no progress made in solving this problem until 1938. Then

Gödel proved that the continuum hypothesis is consistent with the standard axioms of set theory. We are in a subtle area here, where there are *three* possibilities for a theorem: It may be provable, disprovable, or undecidable. Gödel ruled out the “disprovable” possibility.

Kurt Gödel (1906–) is the outstanding mathematical logician of the century. His major achievements include the proof of the completeness of quantification logic, the proof of the incompleteness of standard mathematics, and the proof of the consistency of both the continuum hypothesis and the axiom of choice. The last was announced in *Proc. Nat. Acad. Sci.* (1939) with a full account in *The Consistency of the Continuum Hypothesis*, *Annals of Math. Studies* no. 3, Princeton, 1940.

In 1963 Cohen completed the job by ruling out the “provable” case. So it stands that the continuum problem (it seems better today to say “problem” than “hypothesis”) is undecidable on the basis of the current axioms for set theory.

Paul J. Cohen (1934–), a Chicago Ph.D. of 1958, began his career in analysis. His work on the continuum hypothesis was announced in the *Proc. Nat. Acad. Sci.* 1963–64 with a complete account in *Set Theory and the Continuum Hypothesis*, Benjamin, 1966.

Must reading for anyone interested in the continuum problem is Gödel’s article *What is Cantor’s continuum problem?* in *Amer. Math. Monthly* 54 (1947), 515–525. (A revised and expanded version appears in pages 258–273 of *Philosophy of Mathematics*, edited by P. Benacerraf and H. Putnam, Prentice-Hall, 1964.) We quote two extracts.

. . . it seems . . . that a complete solution of these problems can be obtained only by a more profound analysis (than mathematics is accustomed to give) of the meanings of . . . “set,” “one-to-one correspondence” . . .

Also, it is very suspicious that, as against the numerous plausible propositions which imply the negation of the continuum hypothesis, not one plausible proposition is known which would imply the continuum hypothesis. Therefore one may in good reason suspect that the role of the continuum problem in set theory will be this, that it will finally lead to the discovery of new axioms which will make it possible to disprove Cantor’s Conjecture.

Somewhat similar views are put forward by Cohen on page 151 of the book cited above. Also recommended is the survey “The Continuum Hypothesis” by R. Smullyan in pages 252–260 of *The Mathematical Sciences* (a collection of essays edited by COSRIMS with the collaboration of G. A. W. Boehm, MIT Press, 1969).

So some day our descendants may decide which of the \aleph 's is equal to c . Perhaps \aleph_1 , \aleph_2 , \aleph_{17} ? At present, none of these is ruled out. There is however one known theorem, due to J. König, which implies, among other things, that c cannot be \aleph_ω . This makes a very good exercise—see Exercise 3.

The *generalized continuum hypothesis* asserts that for any infinite cardinal number d , there is no cardinal number properly between d and 2^d . In terms of the \aleph 's, it states that $2^{\aleph_\lambda} = \aleph_{\lambda+1}$ for every ordinal λ . Its status is similar to that of the ordinary continuum hypothesis.

One thing is certain: The assumption of the generalized continuum hypothesis greatly facilitates the computation of cardinal exponentiation—see Exercises 4 and 5.

EXERCISES

- 1.* Infinite sums and infinite products of cardinals are defined as follows: Let d_i be cardinals, and pick D_i with $\mathfrak{o}(D_i) = d_i$. (For Σd_i , pick the D_i 's disjoint.) Then Σd_i is the cardinal of $\cup D_i$, and Πd_i is the cardinal of the Cartesian product of the D_i 's.

Let d_i, e_i be cardinals with $d_i < e_i$ for all i . Prove that $\Sigma d_i < \Pi e_i$.

(Observe that if every $d_i = 1$ and every $e_i = 2$ we get Theorem 6.)

- 2.* Prove that $\aleph_\omega^{\aleph_0} > \aleph_\omega$. (*Hint*: By Exercise 1 we have

$$\aleph_0 + \aleph_1 + \cdots + \aleph_n + \cdots < \aleph_0 \aleph_1 \cdots \aleph_n \cdots$$

The left side is \aleph_ω ; the right side $\leq \aleph_\omega^{\aleph_0}$.)

- 3.* Prove that $c \neq \aleph_\omega$. (*Hint*: Use Exercise 2.)
- 4.* Let d, e be infinite cardinals with $d > 2^e$ and d not a limit cardinal. Assume the generalized continuum hypothesis. Prove that $d^e = d$. (Compare Exercise 6 in Section 2.6.)
- 5.* Assuming the generalized continuum hypothesis, prove that $\aleph_\omega^{\aleph_0} = \aleph_{\omega+1}$.

4

Basic Properties of Metric Spaces

4.1 DEFINITIONS AND EXAMPLES

The very definition of a metric space requires a previous acquaintance with the real numbers. In this book we are taking the real numbers for granted (we have already referred to them many times in the preceding chapters). At this point we shall merely call attention to a property of the real numbers which will play a critical role repeatedly: The real numbers form a chain in which any set bounded from above has a least upper bound, and dually, any set bounded from below has a greatest lower bound. (This property came up in Exercise 11 of Section 1.3 and several later exercises, and was called “conditional completeness.”) A systematic notation for the least upper bound is desirable and, as is customary in analysis, we shall use *sup*. Thus, if A is a set of real numbers bounded from above, we write $\sup A$ for the least upper bound of A . Occasionally we shall also use the notation $\sup A$ when A is not bounded from above; then we write $\sup A = \infty$. If the elements of A are a_i , i ranging over an index set I , we may write $\sup a_i$ for $\sup A$, and we may or may not decorate the symbol $\sup a_i$ with $i \in I$. For greatest lower bound we similarly use *inf*. If A is finite we may use *max* and *min* in place of *sup* and *inf*.

One of the main motivations for studying metric spaces is the desire

to abstract the familiar properties of distance in Euclidean space. We actually assume remarkably little in the definition of a metric space: three quite trivial axioms, and the proposition of Euclid asserting that the sum of two sides of a triangle exceeds the third.

DEFINITION. A metric space is a set M equipped with a real-valued function $D(a, b)$ defined for all $a, b \in M$ so as to satisfy:

- I. $D(a, a) = 0$ for all a .
- II. $D(a, b) > 0$ for $a \neq b$.
- III. $D(a, b) = D(b, a)$.
- IV. $D(a, c) \leq D(a, b) + D(b, c)$ (triangle inequality).

Remarks: 1. To emphasize the switch from pure set theory to a geometrical setup, we usually call the elements of a metric space *points*.

2. To avoid the confusion of too many parentheses, we shall sometimes use brackets for the distance function: $D[\ , \]$.

We proceed to give a small number of examples. A more thorough discussion of examples has been placed in Appendix 1, and should be consulted as needed.

Example 1: The basic example is the set of all real numbers, with $D(a, b) = |a - b|$.

Example 2: A blanket comment applies to all examples, and indeed to all metric spaces. If M is a metric space, and S a subset of M , we may use in S the very same distance function, except that it is restricted to S . Obviously S becomes in this way a metric space. In particular, any set of real numbers is a metric space relative to the distance function $D(a, b) = |a - b|$.

Example 3: Let M be the usual Euclidean plane: the set of all ordered pairs of real numbers. Two typical points of M are $u = (x_1, y_1)$, $v = (x_2, y_2)$. The Euclidean distance is given by

$$(a) \quad D(u, v) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}.$$

Here are two somewhat different possibilities for a distance function in M :

$$(b) \quad D(u, v) = \max(|x_1 - x_2|, |y_1 - y_2|),$$

$$(c) \quad D(u, v) = |x_1 - x_2| + |y_1 - y_2|.$$

Some idea of what these distance functions are like is obtained by plotting in the plane the points at distance 1 from the origin. The result is shown in Figure 10.

Example 4: The final example in this section is an arbitrary set equipped with a trivial distance function. If M is any set, take $D(a, a) = 0$ and $D(a, b) = 1$ for $a \neq b$ in M .

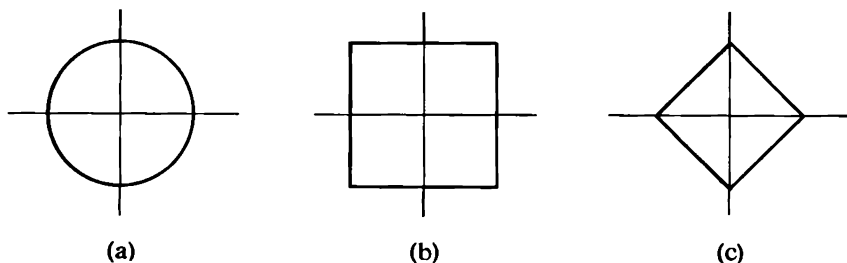


Figure 10

Theorem 26 is a variant of the triangle inequality. Like the triangle inequality, it appears in Euclid's *Elements*.

THEOREM 26. For any points a, b, c in a metric space we have

$$(15) \quad |D(a, c) - D(b, c)| \leq D(a, b).$$

Proof: Because of the symmetry in (15) between a and b , we can assume $D(a, c) \geq D(b, c)$. Then $D(a, c) - D(b, c) \geq 0$, so that we can remove the absolute values in (15). We then find that (15) coincides with the triangle inequality IV.

Two definitions conclude this section.

DEFINITION. The *diameter* of a metric space M is $\sup D(a, b)$, taken over all $a, b \in M$. The diameter may be infinite. If it is finite we of course say that M has *finite diameter*. We write $\text{diam}(M)$ as a shorthand for the diameter of M . When we speak of the diameter of a subset A of M we of course mean $\text{diam}(A)$ in the metric A inherits from M .

In later sections we shall discuss many properties that a mapping between metric spaces may or may not have. At this point we introduce the simplest of these properties.

DEFINITION. A mapping between two metric spaces is an *isometry* if it preserves distance. In symbols: $f: X \rightarrow Y$ is an isometry if $D[f(a), f(b)] = D(a, b)$ for all $a, b \in X$.

Remark: This is the first time we have discussed two metric spaces at once. To avoid unnecessary complications of notation we use the same symbol D in each space, as long as no confusion is imminent.

EXERCISES

1. Let M be a set with two members: a and b . Define $D(a, a) = D(b, b) = 0$, $D(a, b) = D(b, a) = r$, where r is a positive real number. Prove that M is a metric space relative to the function D .
2. Let M be a set with three elements: a, b , and c . Define D so that $D(x, x) = 0$ for all x , $D(x, y) = D(y, x) = a$ positive real number for $x \neq y$. Say $D(a, b) = r$, $D(a, c) = s$, $D(b, c) = t$ with $r \leq s \leq t$. Prove that D makes M a metric space if and only if $t \leq r + s$.
3. On an arbitrary set M suppose that $D(x, x) = 0$, and for $x \neq y$, $D(x, y) = D(y, x)$ is a number between 1 and 2 (the end points 1 and 2 are permitted). Prove that M is a metric space.
4. Let M be a set with a real function D satisfying $D(a, a) = 0$, $D(a, b) \neq 0$ for $a \neq b$, $D(a, b) + D(b, c) \geq D(c, a)$. Prove that D makes M a metric space. (Note: $D(a, b) \geq 0$ and $D(a, b) = D(b, a)$ are not being assumed; they must be proved.)
5. In any metric space prove that $D(a_1, a_n) \leq D(a_1, a_2) + D(a_2, a_3) + \cdots + D(a_{n-1}, a_n)$.
6. A metric space M is called *ultra-metric* if $D(a, c) \leq \max [D(a, b), D(b, c)]$ for all $a, b, c \in M$. Prove that Example 4 is ultrametric. Is the real line ultrametric?
- 7.* Let $f(x)$ be a real-valued function of a real variable, defined for all $x \geq 0$. f is said to be *concave* if for any nonnegative real numbers a, b with $a + b = 1$ we have

$$f(ax + by) \geq af(x) + bf(y).$$

Suppose that f is concave, that $f(0) = 0$, that $f(x) > 0$ for $x > 0$, and that f is monotone in the weak sense, i.e., $x \leq y$ implies $f(x) \leq f(y)$. Let M be a metric space with distance function D . Prove that $f(D)$ is also a distance function on M .

- 8.* If f (a real function of a real variable) has a second derivative satisfying $f'' \leq 0$, prove that f is concave.
- 9.* By using Exercises 7 and 8, or otherwise, prove the following: If D is a metric so is $f(D)$ for $f(x) = \sqrt{x}$, $f(x) = x/(1+x)$, and $f(x) = \min(x, 1)$.
10. If D is a metric, is D^2 necessarily one?
11. Let D_1 and D_2 be metrics on a single space M . Which of the following are metrics on M : $D_1 + D_2$, $\max(D_1, D_2)$, $\min(D_1, D_2)$?
12. Let M be the metric space of all real numbers, and let $x_0 \in M$. Prove that there exist exactly two isometries of M that leave x_0 fixed.
- 13.* (a) Let the metric space M be a subset of the real line (with the inherited metric). Let a, b, c be any points in M . Prove that of the three numbers $D(a, b)$, $D(a, c)$, $D(b, c)$, the largest is the sum of the other two.

- (b) If M has at least 5 points, prove the converse of part (a). In detail: Let M be a metric space with the property that for any $a, b, c \in M$ the largest of $D(a, b)$, $D(a, c)$, $D(b, c)$ is the sum of the other two. Prove that M is isometric to a subset of the real line, if $\alpha(M) \geq 5$. Give an example of failure when M has 4 points.
14. Let M be a metric space in which the distance function assumes only the values 0, 1, 3. Define $x \sim y$ to mean $D(x, y) \leq 1$. Prove that \sim is an equivalence relation. (Observe also that \sim determines the metric.)
- 15.* Let A be any set of positive real numbers. Prove that there exists a metric space whose nonzero distances constitute exactly the set A .
16. ("Antimetric spaces.") Let D on a set M satisfy the axioms for a metric space except that the triangle inequality is reversed:

$$D(a, c) \geq D(a, b) + D(b, c).$$

Prove that M has at most one point.

17. Give an example of a metric space which admits an isometry with a proper subset of itself. (*Hint*: Try Example 4.)
- 18.* Let R be the metric space of all real numbers. Prove that any isometry between two subsets of R can be extended to an isometry of R onto itself.
19. Let M be a set with a distance function D satisfying the postulates for a metric space, except that axiom II (p. 68) is weakened to $D(a, b) \geq 0$ for all a, b . Define $a \sim b$ to mean $D(a, b) = 0$. Prove that \sim is an equivalence relation. Make the set of equivalence classes into a metric space in a natural way.
20. (a) Suppose that a metric space A is a union $A = B \cup C$ of two subsets of finite diameter. Prove A has finite diameter.
 (b) Generalize part (a) to the union of a finite number of subsets.
 (c) Suppose that, in part (a), $B \cap C$ is nonempty. Prove that $\text{diam}(A) \leq \text{diam}(B) + \text{diam}(C)$.

4.2 OPEN SETS

In Sections 4.2–4.4 we shall be extending to arbitrary metric spaces concepts which are probably familiar to the reader from his previous work in analysis. The first of these concepts is that of an open set. This is based in turn on the idea of an open ball.

DEFINITION. Let M be a metric space, x a point in M , and r a positive real number. The set of all y in M with $D(x, y) < r$ is called the *open ball* with center x and radius r , and is denoted by $S_r(x)$. The set of all y in M with $D(x, y) \leq r$ is called the *closed ball* with center x and radius r . We shall not introduce a notation for the closed ball.

The idea behind the definition of an open set U is that every element of U should be nicely surrounded by elements of U .

DEFINITION. Let U be a subset of a metric space M . We say that U is *open* in M if for every $x \in U$ an entire open ball $S_r(x)$ is contained in U .

Remarks: 1. The definition implies (vacuously as it were) that the empty set is open.

2. The entire space M is an open set.

3. The terminology must not be allowed to fool us. We have to prove that an open ball is open. The idea for the proof is that a sufficiently small ball with center y will lie in $S_r(x)$. The proof is illustrated in Figure 11; we present it formally in Theorem 27.

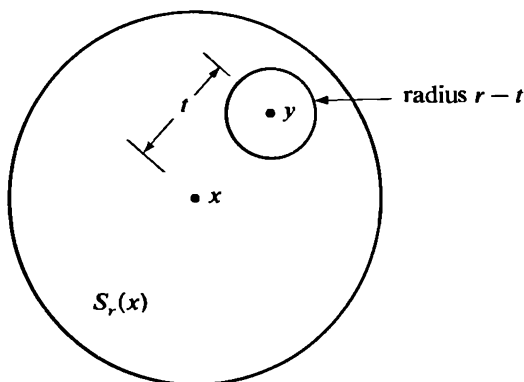


Figure 11

THEOREM 27. Any open ball in a metric space is an open set.

Proof: Let the given open ball be $S_r(x)$. Given $y \in S_r(x)$, we must show that a suitable open ball with center y is contained in $S_r(x)$. Let $D(x, y) = t$; note that $t < r$. Our choice for the radius of the desired ball is $r - t$. If $D(z, y) < r - t$, we have

$$D(z, x) \leq D(z, y) + D(y, x) < t + (r - t) = r.$$

Thus $S_{r-t}(y) \subset S_r(x)$.

The next two theorems describe the behavior of open sets relative to union and intersection.

THEOREM 28. In a metric space any union of open sets is open.

Proof: Let U be the union of the open sets U_i ; we must prove that U is open. Thus given $x \in U$ we must show that a suitable open ball with center x is entirely contained in U . Now x got into U on the authority of being a member of some U_i . In that U_i we have an open ball with center x entirely contained in U_i , and it is contained in U *a fortiori*.

THEOREM 29. *In a metric space a finite intersection of open sets is open.*

Proof: Let U_1, \dots, U_n be open. We are to prove that $U = U_1 \cap \dots \cap U_n$ is open. Let $x \in U$. Then for each $i, x \in U_i$. Since U_i is open we have $S_{r_i}(x) \subset U_i$ for a suitable positive r_i . If we take $r = \min(r_1, \dots, r_n)$, we have $S_r(x) \subset U$.

In Theorem 29 we cannot allow infinite intersections of open sets—see Exercise 3.

We note a quick corollary of Theorems 27 and 28.

THEOREM 30. *A subset of a metric space is open if and only if it is expressible as a union of open balls.*

Proof: Any open ball is an open set (Theorem 27), so by Theorem 28 a union of open balls is open. Conversely, let U be open. For any x in U we have that a suitable $S_r(x)$ is contained in U . Evidently U is the union of all these $S_r(x)$'s.

Remark: In Theorem 30 we should not suppose that we are getting a *disjoint* union of open balls. This is rarely the case. However, it is essentially true for the metric space of real numbers—see Exercise 11.

We introduce a variant of the notion of open set. It is useful in many contexts.

DEFINITION. A *neighborhood* of a point x in a metric space is a subset containing an open set containing x .

Remarks: 1. It is obviously equivalent to require that the neighborhood contain an open ball containing x , or that it contain an open ball with center x .

2. In the terminology of some authors, a neighborhood is required to be open. We shall use the designation *open neighborhood* when we require openness.

3. A useful example of a neighborhood of x is a *closed* ball with center x . The virtually trivial proof of the next theorem is left to the reader.

THEOREM 31. *A subset of a metric space is open if and only if it contains, along with any point x , some neighborhood of x .*

We conclude this section with two definitions.

DEFINITION. A point x in a metric space is called *isolated* if the set $\{x\}$ consisting of x alone is open.

DEFINITION. A metric space is *discrete* if every subset is open. (By Theorem 28 it is equivalent to say that every point is isolated.)

EXERCISES

- (a) Prove that in a metric space the complement of a point is open.
(b) Prove that any set in a metric space is an intersection of open sets.
- Let x and y be distinct points of a metric space M . Prove that there exist in M disjoint open sets U and V with $x \in U$, $y \in V$.
- If a metric space M has the property that every intersection of open sets is open, prove that M is discrete.
- Prove that Example 4 in Section 4.1 is discrete.
- Prove that in the metric space of all real numbers there are no isolated points.
- Let A be a finite open subset of a metric space M . Prove that every point in A is an isolated point of M .
- Given a point x in a metric space M , and positive numbers r, s with $r < s$, prove that the set of all $y \in M$ satisfying $r < D(x, y) < s$ is open. (A reasonable name for such a set is "open ring" or "open annulus.")
- Let x be a point of a metric space M . Prove that the following two statements are equivalent:
 - x is not isolated,
 - Every neighborhood of x contains an infinite number of points of M .
- * Let M be an infinite metric space. Prove that M contains an open set U such that both U and its complement are infinite.
- Let A and B be open balls in an ultrametric space (Exercise 6, Section 4.1). Prove that if A and B have a nonvoid intersection, then one of them contains the other.
- * Let R be the metric space of all real numbers. Prove that any bounded open set in R is a countable union of disjoint open intervals.
- Prove that any nonvoid open subset of the real line has cardinal number c .

4.3 CONVERGENCE; CLOSED SETS

Convergence of a sequence in a metric space is defined just as in elementary analysis.

DEFINITION. Let x_1, x_2, x_3, \dots be a sequence in a metric space M and let $x \in M$. The sequence is said to *converge* to x if for any positive real number ϵ there exists an integer N (depending on ϵ) such that $D(x_i, x) < \epsilon$ for all $i \geq N$. x is called the *limit* of the sequence.

The elements of the sequence need not be distinct, and they need not differ from x . For example, the following three sequences of real numbers all converge to 0:

$$\begin{aligned} &1, \frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \dots \\ &1, 0, \frac{1}{2}, 0, \frac{1}{3}, 0, \frac{1}{4}, 0, \dots \\ &0, 0, 0, 0, \dots \end{aligned}$$

We write $\{x_i\}$ as a handy symbol for the sequence x_1, x_2, x_3, \dots . Strictly speaking, the symbol $\{x_i\}$ should denote the unordered set of elements, but there should be no confusion. We write $x_i \rightarrow x$ for the statement that the sequence $\{x_i\}$ converges to x . This of course has nothing to do with the use of an arrow in describing a function $f: A \rightarrow B$.

The definition can be rephrased so as to avoid mentioning a number ϵ by saying that a sequence converges to x if every neighborhood of x contains the sequence after deletion of a finite initial segment.

We define closed sets in terms of convergence.

DEFINITION. A set F in a metric space M is said to be *closed* if whenever a sequence $\{x_i\}$ of elements of F converges to $x \in M$, then $x \in F$. Without symbols: A closed set is one which contains all limits of its convergent sequences.

Remarks: M is a closed set. The empty set is a closed set. In fact, M and the empty set are both open and closed. There may be still other sets which are open and closed; this depends upon M . If there are none, M is called *connected* (see Appendix 3). In a discrete space all subsets are open and closed. The real line is connected.

In Theorem 33 we shall obtain a decisive connection between closed sets and open sets. Before doing so, we characterize convergence in terms of open balls.

THEOREM 32. *Let A be any subset and x any point in a metric space. There exists a sequence of elements of A converging to x if and only if every $S_r(x) \cap A$ is nonempty.*

Proof: If the sequence x_1, x_2, x_3, \dots of elements of A converges to x , then, for any given r , $x_i \in S_r(x)$ for sufficiently large i . Conversely, suppose that each $S_r(x) \cap A$ is nonempty. Then, in particular, $S_{1/n}(x) \cap A$ is nonempty for every positive integer n . Pick $x_n \in S_{1/n}(x) \cap A$. We claim that the sequence x_1, x_2, x_3, \dots converges to x . Indeed, given $\epsilon > 0$ we need only pick N large enough so that $1/N < \epsilon$ in order to be sure that $D(x_i, x) < \epsilon$ for all $i \geq N$.

Remarks: 1. The countable axiom of choice slipped into the proof of Theorem 32 without much fanfare.

2. Theorem 32 can be restated in terms of neighborhoods, as in Theorem 36.

THEOREM 33. *A subset of a metric space is closed if and only if its complement is open.*

Proof: Let F be a closed set in the metric space M , and let U be the complement of F in M . We have to show that U is open. That is, given $x \in U$, we have to show that some $S_r(x)$ is contained in U . If this is false then every $S_r(x) \cap F$ is nonempty. By Theorem 32, this entails the existence of a sequence of elements of F converging to x . Since F is closed, we get the contradiction $x \in F$.

Conversely, suppose F is the complement of an open set U . To prove that F is closed we take a sequence x_1, x_2, x_3, \dots of elements of F converging to x and have to prove that $x \in F$. Suppose on the contrary that $x \in U$. Then some $S_r(x) \subset U$. For large enough i , $D(x_i, x) < r$, $x_i \in S_r(x)$, a contradiction.

Using Theorem 33 in conjunction with Theorems 28 and 29, we obtain the next two theorems simply by set-theoretic complementation.

THEOREM 34. *The intersection of any collection of closed sets in a metric space is again a closed set.*

THEOREM 35. *The union of a finite number of closed sets in a metric space is again a closed set.*

The reader might find it enlightening to prove Theorems 34 and 35 directly from the definition of a closed set.

The fact that the intersection of any number of closed sets is closed makes possible a certain general construction. Let A be any subset of a metric space M . There exist closed sets containing A , for instance M . Let \bar{A} denote the intersection of all closed subsets of M containing A . Then \bar{A} is closed, and in an obvious sense it is the smallest closed set containing A . We call \bar{A} the *closure* of A .

We can identify the closure of a set by convergent sequences or by neighborhoods.

THEOREM 36. *Let A be a subset and x a point in a metric space M . The following three statements are equivalent:*

- (a) x lies in the closure \bar{A} of A .
- (b) Every neighborhood of x intersects A .
- (c) There exists a sequence of elements of A converging to x .

Proof: The equivalence of (b) and (c) is essentially just a rephrasing of Theorem 32. We shall therefore confine ourselves to proving (a) \Rightarrow (b) and (c) \Rightarrow (a).

(a) \Rightarrow (b): We are given $x \in \bar{A}$ and wish to show that any neighborhood U of x intersects A . Since we can shrink U to an open neighborhood of x , we might as well assume that U is open. Suppose that $A \cap U$ is empty. If T is the complement of U we have that T is closed (Theorem 33) and $A \subset T$. Then $\bar{A} \subset T$, $x \in T$, a contradiction.

(c) \Rightarrow (a): If a sequence of elements of A converges to x , then x lies in any closed set containing A , $x \in \bar{A}$.

We introduce an additional concept, which differs slightly from asserting that a point lies in the closure of a set.

DEFINITION. Let A be a subset and x a point in a metric space. We say that x is a *limit point* of A if it lies in the closure of $A - \{x\}$.

In greater detail: If $x \notin A$, then to say that x is a limit point of A is the same as saying that it lies in the closure \bar{A} of A . If $x \in A$, then (although x of course lies in \bar{A}) we call x a limit point of A only if it lies in the closure of the set obtained by deleting x from A .

The reason for introducing limit points is that in Section 5.2 we shall find the following statement significant: any infinite set has a limit point. Without the shift of meaning, this statement would always be trivially true.

For a characterization of limit points, see Exercise 5.

Convergence in a metric space can be closely tied to convergence of real numbers. For instance, $x_n \rightarrow x$ if and only if $D(x_n, x) \rightarrow 0$, an easy

observation that we leave the reader as Exercise 2. We shall need Theorem 37, a closely related result.

THEOREM 37. *If $x_i \rightarrow x$ and $y_i \rightarrow y$ in a metric space, then*

$$D(x_i, y_i) \rightarrow D(x, y).$$

Proof: We have

$$D(x_i, y_i) \leq D(x_i, x) + D(x, y) + D(y, y_i).$$

Hence, given any $\epsilon > 0$, we have

$$D(x_i, y_i) < D(x, y) + 2\epsilon$$

for sufficiently large i . Similarly,

$$D(x, y) < D(x_i, y_i) + 2\epsilon$$

for sufficiently large i . Combining these inequalities, we obtain

$$|D(x_i, y_i) - D(x, y)| < 2\epsilon.$$

Hence $D(x_i, y_i) \rightarrow D(x, y)$.

EXERCISES

1. If in a metric space a sequence converges to both x and y , prove that $x = y$.
2. Prove: In a metric space, $x_i \rightarrow x$ if and only if $D(x_i, x) \rightarrow 0$.
3. Let A and B be subsets of a metric space. Prove: (a) $A \subset B$ implies $\bar{A} \subset \bar{B}$, (b) $\overline{A \cup B} = \bar{A} \cup \bar{B}$. Does $\overline{A \cap B} = \bar{A} \cap \bar{B}$ always hold?
4. Given distinct points x and y in a metric space, prove that there exist open sets U and V such that $x \in U$, $y \in V$, and $\bar{U} \cap \bar{V}$ is empty. (Observe that this sharpens Exercise 2 in Section 4.2.)
5. Let A be a subset and x a point in a metric space. Prove that x is a limit point of A if and only if there exists a sequence $\{x_i\}$ of elements of A converging to x , with all $x_i \neq x$. Prove further that the x_i 's can be selected so that the numbers $D(x_i, x)$ decrease strictly monotonically, i.e.,

$$D(x_1, x) > D(x_2, x) > D(x_3, x) > \dots$$

6. (a) Prove that a closed ball in a metric space is a closed set.
 (b) The *sphere* with radius r and center x is the set of all y with $D(x, y) = r$. Prove that a sphere is a closed set.
 (c) Given a point x in a metric space M , and positive numbers r, s with $r < s$, prove that the set of all $y \in M$ satisfying $r \leq D(x, y) \leq s$ is closed. (Compare with Exercise 7 in Section 4.2. This is a "closed annulus.")

7. Prove that in a metric space a set consisting of just one point is closed. (*Remark: Do not get this by complementation from Exercise 1(a) in Section 4.2!*)
8. Prove: The diameter of a set A in a metric space is the same as the diameter of its closure.
9. Let A be a subset of a metric space M .
 - (a) If U is an open set in M , prove that $U \cap A$ is open in the metric space A .
 - (b) If F is a closed set in M , prove that $F \cap A$ is closed in A .
 - (c) Prove that any open (resp. closed) set in A is obtainable in this way from an open (resp. closed) set in M .
10. (Distance between sets): Let A and B be sets in a metric space. Define $D(A, B)$ to be $\inf D(a, b)$ where the \inf is taken over all $a \in A, b \in B$. Observe that $D(A, B) = 0$ if A and B have a point in common, but that this condition is not necessary.
 - (a) If B consists of a single point x , prove that $D(A, B) = 0$ if and only if x is in the closure of A .
 - (b) Give an example where A and B are both closed, $A \cap B$ is empty, and $D(A, B) = 0$. (*Suggestion: a hyperbola and its asymptotes.*)
- 11.* Prove that in a metric space the closure of a countable set has cardinal number at most c ($= 2^{\aleph_0}$, the cardinal number of the continuum).
- 12.* Prove that the following statements are equivalent for a metric space M :
 - (a) Every subset of M is either open or closed.
 - (b) At most one point of M is not isolated.
- 13.* Let M be a metric space in which the closure of every open set is open. Prove that M is discrete.
- 14.* Prove that in a metric space every open set is the union of a countable number of closed sets, and every closed set is the intersection of a countable number of open sets.
- 15.* Prove that any infinite metric space contains an infinite subset which is discrete (as a metric space in its own right).
16. Prove that a metric space is discrete if and only if every convergent sequence is ultimately constant.
17. Prove that a metric space is discrete if and only if it has no limit points.
18. Let $\{x_i\}$ be a sequence of distinct elements in a metric space, and suppose that $x_i \rightarrow x$. Let f be a one-to-one map of the set of x_i 's into itself. Prove that $f(x_i) \rightarrow x$.
- 19.* If a metric space M has only countably many open sets, prove that M is countable.
- 20.* Prove that \aleph_1 disjoint closed sets in a metric space can be enlarged to disjoint open sets.

4.4 CONTINUITY

The definition of continuity is the familiar epsilon-delta one of calculus.

DEFINITION. Let X and Y be metric spaces, let f be a function from X to Y , let x_0 be a point in X , and write $y_0 = f(x_0)$. We say that f is continuous at x_0 if the following is true: For any $\epsilon > 0$ there exists $\delta > 0$ such that $D(x, x_0) < \delta$ implies $D[f(x), y_0] < \epsilon$.

We can characterize continuity at a point in terms of neighborhoods, or in terms of convergence. The characterization in terms of neighborhoods is just a slight change of language.

THEOREM 38. *Let X and Y be metric spaces, let f be a function from X to Y , let x_0 be a point in X , and write $y_0 = f(x_0)$. Then f is continuous at x_0 if and only if the following is true: For any neighborhood V of y_0 there exists a neighborhood U of x_0 with $f(U) \subset V$.*

Proof: Suppose that f is continuous at x_0 . The neighborhood V of y_0 contains an open ball with center y_0 , say of radius ϵ . For the δ which corresponds to this ϵ (in the definition of continuity at x_0) we have $f(S_\delta(x_0)) \subset S_\epsilon(y_0)$. We can therefore take $U = S_\delta(x_0)$.

Suppose, conversely, that the condition in the final phrase of the theorem is satisfied. We have to prove that f is continuous at x_0 . Challenged with an ϵ , we have to produce a suitable δ . Take $V = S_\epsilon(y_0)$. The corresponding neighborhood U of x_0 contains some ball $S_\delta(x_0)$. This δ will do.

The characterization of continuity by convergence of sequences is not quite such a formality. In half the proof the countable axiom of choice sneaks in.

THEOREM 39. *Let X, Y, f, x_0 , and y_0 be as in Theorem 38. Then f is continuous at x_0 if and only if the following is true: Whenever a sequence $\{x_i\}$ of elements of X converges to x_0 , and $y_i = f(x_i)$, the sequence $\{y_i\}$ converges to y_0 .*

Proof: Suppose that f is continuous at x_0 and let x_i converging to x_0 be given. For large i , $D(x_i, x_0)$ is arbitrarily small, and hence (by the assumed continuity) $D(y_i, y_0)$ is arbitrarily small for large i . Therefore $y_i \rightarrow y_0$.

Conversely, we assume the preservation of convergence, and seek to prove that f is continuous at x_0 . The proof is indirect. If f is not continuous

at x_0 , then some ϵ cannot be matched by a δ . In particular the choice $\delta = 1/n$ will not do, for any positive integer n . In detail, we find this to mean that there exists x_n with $D(x_n, x_0) < 1/n$ and $D(y_n, y_0) > \epsilon$, where $y_n = f(x_n)$. Then the sequence $\{x_i\}$ converges to x_0 while the sequence $\{y_i\}$ does not converge to y_0 , a contradiction.

We next define continuity "in the large."

DEFINITION. A function f from a metric space X to a metric space Y is *continuous* if it is continuous at every point of X .

Global continuity has a remarkably simple formulation in terms of open or closed sets.

THEOREM 40. *Let f be a mapping from a metric space X to a metric space Y . Then the following three statements are equivalent:*

- (a) f is continuous.
- (b) The complete inverse image of an open set is open, i.e., for any open set V in Y , $f^{-1}(V)$ is an open subset of X .
- (c) The complete inverse image of a closed set is closed, i.e., for any closed set G in Y , $f^{-1}(G)$ is a closed subset of X .

Proof: Since (Theorem 33) open and closed sets are complements of each other, the equivalence of (b) and (c) is an easy matter of pure set theory. We shall therefore confine ourselves to proving the equivalence of (a) and (b).

(a) \Rightarrow (b): Let V be open in Y and let x be a point in $f^{-1}(V)$. Then $y = f(x)$ lies in V . Since V is a neighborhood of y , we have $f(U) \subset V$ for a suitable neighborhood U of x (Theorem 38). Since $U \subset f^{-1}(V)$, the criterion of Theorem 31 shows that $f^{-1}(V)$ is open.

(b) \Rightarrow (a): To prove that f is continuous we prove continuity at x . With $y = f(x)$, we take an open neighborhood V of y . Then $f^{-1}(V)$ is open by hypothesis and therefore contains a neighborhood U of x (we may even take $U = f^{-1}(V)$). Since $f(U) \subset V$ we have fulfilled the criterion of Theorem 38.

In Section 5.1 we shall consider the question of extending a continuous function from a set on which it is defined to a larger set. We shall conclude this section with a theorem on the uniqueness of such an extension. The following definition is relevant.

DEFINITION. A subset A of a metric space M is *dense* in M if $\bar{A} = M$

THEOREM 41. *Let f and g be continuous functions from a metric space X to a metric space Y . Let B be the set of all x in X for which $f(x) = g(x)$. Then B is a closed subset of X .*

Proof: Suppose the elements b_i lie in B and converge to c . We have to prove that $c \in B$. Now by Theorem 40, $f(b_i) \rightarrow f(c)$ and $g(b_i) \rightarrow g(c)$. Since $f(b_i) = g(b_i)$ we deduce $f(c) = g(c)$ (there is a small point here which is covered by Exercise 1 in Section 4.3). Hence $c \in B$ and B is closed.

THEOREM 42. *Let A be a dense subset of a metric space X , and let f and g be continuous functions from X to a metric space Y . Suppose that f and g coincide on A . Then they coincide on all of X .*

Proof: The set where f and g agree is closed by Theorem 41, contains A , and therefore contains its closure X .

EXERCISES

1. Let f be a function from a metric space to a second metric space. Prove that f is continuous if and only if it sends convergent sequences into convergent sequences.
2. Let u be a fixed point of a metric space M . The function $f(x) = D(u, x)$ maps M into the real numbers. Prove that f is continuous.
3. Let A be a fixed subset of a metric space M . The function $f(x) = D(A, x)$ maps M into the real numbers—see Exercise 10 in Section 4.3. Prove that f is continuous.
4. Let A be a closed subset and y a point in a metric space M , $y \notin A$. Prove that there exists a continuous real function on M which vanishes on A but not at y . (Use the preceding exercise and Exercise 10(a) in Section 4.3.)
5. Let X and Y be metric spaces, f a function from X to Y .
 - (a) If X is a union of open sets U_i on each of which f is continuous, prove that f is continuous on X .
 - (b) If X is a finite union of closed sets F_1, \dots, F_n on each of which f is continuous, prove that f is continuous on X .
 - (c) Can part (b) be extended to an infinite number of closed sets?
6. Prove that a metric space X is discrete if and only if every function on X (to an arbitrary second metric space) is continuous.
7. (a) Find the cardinal number of the set of all continuous functions from the real numbers to itself. (*Hint:* Any such function is determined by its values at the rational numbers.)

- (b) Find the cardinal number of the set of all continuous functions from the Euclidean plane to itself.
8. If $f: X \rightarrow Y$ is continuous at x_0 , and $g: Y \rightarrow Z$ is continuous at $f(x_0)$, prove that $gf: X \rightarrow Z$ is continuous at x_0 .
- 9.* Two metric spaces are said to be *homeomorphic* if there exists between them a one-to-one correspondence which is continuous in both directions. Prove that any metric space is homeomorphic to one of finite diameter. (*Hint*: See Exercise 9 in Section 4.1.)
- 10.* (a) A mapping from one metric space to a second is called *open* if it sends open sets into open sets. Give an example of a continuous function which is not open and an open function which is not continuous.
- (b) Similarly a mapping is called *closed* if it sends closed sets into closed sets. Give examples, as in part (a), to distinguish closed mappings from continuous ones.
- (c) Let X and Y be metric spaces, and let f be a function from X to Y which is one-to-one and onto. Prove that the following three statements are equivalent: f is open, f is closed, f^{-1} is continuous.

5

Completeness, Separability, and Compactness

5.1 COMPLETENESS

We next introduce, in the context of quite general metric spaces, the concept of a Cauchy sequence. The idea is to extract the essence of convergence, without actually having an element toward which the convergence is taking place. This is done by postulating that, eventually, the elements of the sequence get arbitrarily close to each other.

DEFINITION. A sequence $\{x_i\}$ in a metric space is called *Cauchy* if the following is true: For any positive ϵ there exists a positive integer N such that $D(x_i, x_j) < \epsilon$ for all $i, j \geq N$.

It is true that any convergent sequence is Cauchy.

THEOREM 43. *Any convergent sequence in a metric space is a Cauchy sequence.*

Proof: Suppose that $x_i \rightarrow y$. Given ϵ , pick N large enough so that $D(x_i, y) < \epsilon/2$ for $i \geq N$. Then for $i, j \geq N$,

$$D(x_i, x_j) \leq D(x_i, y) + D(y, x_j) < \epsilon/2 + \epsilon/2 = \epsilon.$$

As a first step in getting some kind of converse to Theorem 43 we assume that a given Cauchy sequence has a convergent subsequence. Since this is the first time the notion of "subsequence" has come up, we shall take a moment to try to make sure there is no possibility of misunderstanding. For a sequence consisting of distinct elements, ambiguity is highly unlikely. Let us illustrate with an example where there are repeated elements. The sequence

$$1, 0, \frac{1}{2}, 0, \frac{1}{3}, 0, \frac{1}{4}, 0, \dots$$

has as one subsequence

$$1, \frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \dots$$

Another subsequence is

$$0, 0, 0, 0, \dots$$

A more formal description of a subsequence can be given as follows. Look at a sequence of elements in a set A as a function ϕ from the positive integers Z to A . Let h be a one-to-one monotone function from Z to Z . The composite of the functions h and ϕ

$$Z \xrightarrow{h} Z \xrightarrow{\phi} A$$

names a typical subsequence of the sequence given by ϕ .

THEOREM 44. *If a Cauchy sequence $\{x_i\}$ in a metric space has a subsequence convergent to y , then the whole sequence converges to y .*

Proof: Given $\epsilon > 0$, pick N so that $D(x_i, x_j) < \epsilon/2$ for $i, j \geq N$. We can find an element x_k of the subsequence with $k \geq N$ and with k large enough so that $D(x_k, y) < \epsilon/2$. Then for all $i \geq N$

$$D(x_i, y) \leq D(x_i, x_k) + D(x_k, y) < \epsilon/2 + \epsilon/2 = \epsilon.$$

We prove next that any Cauchy sequence is "bounded."

THEOREM 45. *Let x_1, x_2, x_3, \dots be a Cauchy sequence in a metric space. Then the set consisting of the x_i 's has finite diameter.*

Proof: Given $\epsilon > 0$, we have an N such that $D(x_i, x_j) < \epsilon$ for $i, j \geq N$. Thus beyond x_N we have the estimate ϵ for the diameter. The addition of the finite number of elements x_1, \dots, x_{N-1} keeps the diameter finite (compare Exercise 20 in Section 4.1).

A Cauchy sequence need not converge. For instance, in the metric space of rational numbers take a sequence approaching $\sqrt{2}$. This is a Cauchy sequence, but it fails to converge to anything in our given metric

space. True, it “converges” to a ghost-like object in a larger metric space. This suggests two things: We should have a notion of completeness of a metric space, and a way of enlarging any metric space to be complete.

DEFINITION. A metric space M is *complete* if every Cauchy sequence in M is convergent (to a point of M , of course).

The primary example of a complete metric space is the space of real numbers. We shall prove this now, or rather we shall reduce it to a statement which we shall take as known: that any bounded monotone sequence of real numbers converges. Theorem 46 is a preliminary elementary result on chains. (In Theorem 46 we could allow repetitions in the sequence, and then the proof of Theorem 47 would be slightly simpler.)

THEOREM 46. *Let x_1, x_2, x_3, \dots be a sequence of distinct elements in a chain. Then there is either a subsequence which is monotone increasing or a subsequence which is monotone decreasing.*

Proof: Assume that no subsequence is monotone increasing. If we start with any x_i , hunt up an x_j with $j > i$ and $x_i < x_j$, then an x_k with $k > j$ and $x_j < x_k$, etc., this process must end in a finite number of steps. It ends with an element which is larger than all later ones in the sequence. Let us (just for this proof) call such an element *dominant*. What we have just shown is that after any point in the sequence we can find a dominant element. Therefore there exists a subsequence which consists of dominant elements. Necessarily this subsequence is monotone decreasing.

THEOREM 47. *The metric space of real numbers is complete.*

Proof: Let $\{x_i\}$ be a Cauchy sequence of real numbers. We set out to find a convergent subsequence. If some number is repeated infinitely often, there is our subsequence. So we may assume that every number gets repeated only a finite number of times. It follows that there is a subsequence consisting of distinct real numbers. Then by Theorem 46, we can drop down further to a subsequence which is monotone (either increasing or decreasing), and bounded by Theorem 45. As we remarked above, we are taking it as known that such a subsequence converges, say to y . By Theorem 44, the whole sequence converges to y .

From the real line we should go on to Euclidean space, but we shall postpone this till the consideration of Cartesian products (Section 6.1).

Closed subspaces of complete spaces are complete.

THEOREM 48. *Let A be a closed subset of a complete metric space M . Then A is complete.*

Proof: Let $\{a_i\}$ be a Cauchy sequence in A . It can be regarded as a Cauchy sequence in M , and as such it converges to an element b in M . Since A is closed, $b \in A$.

Conversely, complete subspaces are closed (whether or not the big space is complete).

THEOREM 49. *Let A be a subset of a metric space M . Suppose that A is complete (in the induced metric). Then A is a closed subset of M .*

Proof: Given $\{a_i\}$ in A converging to $b \in M$, we have to show that $b \in A$. Now the sequence $\{a_i\}$ is Cauchy (Theorem 43). Since A is complete, $\{a_i\}$ converges to a limit in A . Hence $b \in A$.

We prove a characterization of completeness in metric spaces which generalizes the "nested sequences" often used on the real line.

THEOREM 50. *For a metric space M the following two statements are equivalent:*

- (a) M is complete.
- (b) For any descending sequence $A_1 \supset A_2 \supset A_3 \supset \cdots$ of nonempty closed sets with diameters approaching 0, the intersection $\cap A_i$ is nonempty.

Remarks: 1. In general, if $\text{diam}(A_i) \rightarrow 0$, then $\cap A_i$ is either empty or contains exactly one point.

2. If in (b) we drop the assumption $\text{diam}(A_i) \rightarrow 0$, we get a stronger property which we shall call compactness in Section 5.3. Note in this connection Exercise 11.

Proof: (a) \Rightarrow (b): For each i , pick a point $a_i \in A_i$. Given $\epsilon > 0$ we can find N large enough so that $\text{diam}(A_N) < \epsilon$. Then $D(a_i, a_j) < \epsilon$ for $i, j \geq N$. Hence the sequence $\{a_i\}$ is Cauchy. Since M is complete, $a_i \rightarrow b$ for some $b \in M$. We claim $b \in \cap A_i$. For take a fixed A_k . From a_k on, the points of the sequence lie in A_k . Since A_k is closed, $b \in A_k$.

(b) \rightarrow (a): Let $\{x_i\}$ be a Cauchy sequence in M ; we have to prove that the sequence converges to a point in M . Let B_i be the set consisting of the points x_i, x_{i+1}, \dots . We have $B_1 \supset B_2 \supset B_3 \supset \cdots$; also $\text{diam}(B_i) \rightarrow 0$, as follows readily from the fact that the x 's form a Cauchy sequence. Let A_i be the closure of B_i . The A 's inherit from the B 's the

properties that $A_1 \supset A_2 \supset A_3 \supset \dots$ and that $\text{diam}(A_i) \rightarrow 0$ (note that $\text{diam}(A_i) = \text{diam}(B_i)$ by Exercise 8 in Section 4.3). By hypothesis there is an element y in $\bigcap A_i$. It is evident that the sequence $\{x_i\}$ converges to y .

We turn our attention to the completion of a metric space. We first discuss the uniqueness, and then the existence. For the uniqueness, the concept of uniform continuity is pertinent.

DEFINITION. A function f from a metric space X to a metric space Y is *uniformly continuous* if for any $\epsilon > 0$ there exists a $\delta > 0$ such that, for all $x, x' \in X$, $D(x, x') < \delta$ implies $D[f(x), f(x')] < \epsilon$.

The idea behind this definition is that the epsilon-delta definition of continuity is being applied globally, rather than one point at a time. In other words, for uniform continuity it is insisted that the selected δ work everywhere at once in X . Let us give an example to illustrate the point. As we shall see in Section 5.3, no such example can be given with X a closed interval—or any bounded closed set on the real line. A suitable space is obtained by omitting an end point from a closed interval. It is simpler still to take X to be a sequence of points rather than an interval. Our choice for X is

$$1, \frac{1}{2}, \frac{1}{3}, \dots, \frac{1}{n}, \dots$$

in its usual metric as a subset of the real line. Now X is a discrete space; every point is open, and so every subset is open. Hence any function whatever from X to any metric space is continuous. This simplification enables us to concentrate entirely on the matter of uniform continuity. We define f on X by $f(1/n) = 1$ for n odd and 0 for n even. Then no matter how small we choose δ there are two points in X within δ of each other where f takes the values 0 and 1. Hence $\epsilon = 1$ cannot be matched by any δ . In this connection, see also Exercises 12 and 13.

We begin the theory of uniformly continuous functions by observing their effect on Cauchy sequences.

THEOREM 51. *A uniformly continuous function carries Cauchy sequences into Cauchy sequences.*

Proof: Let the function be f and let $\{x_i\}$ be a Cauchy sequence. We must prove that $\{f(x_i)\}$ is a Cauchy sequence. Given $\epsilon > 0$, produce the δ that is guaranteed by uniform continuity. Then pick N so that

$D(x_i, x_j) < \delta$ for $i, j \geq N$. We deduce that $D[f(x_i), f(x_j)] < \epsilon$ for $i, j \geq N$.

Next comes the extendibility of uniformly continuous functions.

THEOREM 52. *Let A be a dense subset of the metric space B . Let f be a uniformly continuous function from A to a complete metric space C . Then f has a unique extension to a continuous function from B to C . The extended function is uniformly continuous. Furthermore, if f is an isometry on A , the extended function is likewise an isometry.*

Proof: The uniqueness was already proved in Theorem 42; it has nothing to do with uniform continuity or completeness.

Let b be any point in B . Since A is dense in B , there is a sequence $\{a_i\}$ of elements of A converging to b . By Theorem 51, $\{f(a_i)\}$ is a Cauchy sequence in C . Since C is complete, this Cauchy sequence converges to $c \in C$. We plan to extend our function by sending b into c ; we shall use the same letter f for the extended function, so we are proposing $f(b) = c$. It is most urgent to prove that the definition is independent of the choice of the sequence converging to b . We shall prove this by an "interlacing" technique. Let $\{\alpha_i\}$ be another sequence of elements of A converging to b . We observe that the sequence

$$a_1, \alpha_1, a_2, \alpha_2, a_3, \dots$$

converges to b . By Theorem 52 again, the sequence

$$(16) \quad f(a_1), f(\alpha_1), f(a_2), f(\alpha_2), f(a_3), f(\alpha_3), \dots$$

is a Cauchy sequence in C . Since a subsequence of (16) converges to c , the whole sequence converges to c by Theorem 44. We have thus shown that f is well-defined. Of course f , defined this way on B , coincides with the given f on A .

We proceed to prove that the extended function f is uniformly continuous. Given $\epsilon > 0$, we select the δ that works for $f: A \rightarrow C$, not for the given ϵ , but for $\epsilon/2$. We claim that δ is a good choice for checking uniform continuity of the extended f . So: Let b and β be given in B with $D(b, \beta) < \delta$. We assert that $D[f(b), f(\beta)] < \epsilon$. To prove this, we take sequences $\{a_i\}$, $\{\alpha_i\}$ in A converging to b and β . By Theorem 37, $D(a_i, \alpha_i) \rightarrow D(b, \beta)$. Hence $D(a_i, \alpha_i) < \delta$ for large i , and for the same large values of i we have $D[f(a_i), f(\alpha_i)] < \epsilon/2$. From Theorem 37 again we deduce $D[f(b), f(\beta)] \leq \epsilon/2 + \epsilon/2 = \epsilon$. (In the limit the inequality $< \epsilon/2$ can degenerate to $\leq \epsilon/2$. It was solely for this reason that $\epsilon/2$ was introduced at the beginning of this paragraph.) We have proved the uniform continuity of the extended f .

The final statement of the theorem is quickly proved in a similar way. Using the notation above, we have

$$(17) \quad D(b, \beta) = \lim D(a_i, \alpha_i) = \lim D[f(a_i), f(\alpha_i)] = D[f(b), f(\beta)],$$

where the first and last equations in (17) follow from Theorem 37, and the middle equality is a consequence of the assumption that f is an isometry on A . Thus the extended f is an isometry on B .

We define the completion of a metric space.

DEFINITION. Let X and Y be metric spaces. We say that Y is a *completion* of X if Y is complete and contains X as a dense subset.

As we stated above, uniqueness will be proved before existence. Uniqueness follows quickly from Theorem 52.

THEOREM 53. *Let Y and Z be completions of a metric space X . Then there exists an isometry of Y onto Z which is the identity on X .*

Proof: We apply Theorem 52 with f the identity map of X into Z . There results a uniformly continuous map f of Y into Z . By the last part of Theorem 52, f is an isometry on Y . It remains to argue that the range of f is all of Z . Now any point z in Z is the limit of a sequence $\{x_i\}$ in X . If $\{x_i\}$ converges to y in Y , then $f(y) = z$.

Now that uniqueness has been established, we may speak of “the” completion rather than “a” completion.

We shall finish this section with a construction of the completion of a metric space. This can be handled in a head-on fashion, using suitable equivalence classes of Cauchy sequences as points. A thoughtful person might well reflect that this is just the sort of thing one does in constructing the real numbers. *Query:* Having done it once for the real numbers, do we really have to do it again? The answer is “no”; we can instead make use of a device for embedding a metric space into a function space which is known in advance to be complete. While the proof is perhaps not really shorter this way, it is more interesting, and the intermediate steps (Theorems 55 and 56) are valuable for their own sake.

We introduce *uniform convergence*. It is related to pointwise convergence in the same way that uniform continuity is related to continuity: We require an estimate that works everywhere at once.

DEFINITION. Let f_i, f be functions from a metric space X to a metric space Y . We say that the sequence f_i *converges uniformly* to f if the follow-

ing is true: For any $\epsilon > 0$ there exists N (depending on ϵ but independent of x) such that $D[f_i(x), f(x)] < \epsilon$ for all $i \geq N$ and all $x \in X$.

THEOREM 54. *The limit of a uniformly convergent sequence of continuous functions is again continuous.*

Proof: Let f_i, f be functions from X to Y , with f_i continuous and converging uniformly to f . We prove continuity of f at x_0 . Given $\epsilon > 0$ we pick i large enough so that $D[f_i(x), f(x)] < \epsilon/3$ for all x . Then (by the continuity of f_i) pick δ so that $D[f_i(x), f_i(x_0)] < \epsilon/3$ for $D(x, x_0) < \delta$. We combine these estimates to get

$$\begin{aligned} D[f(x), f(x_0)] &\leq D[f(x), f_i(x)] + D[f_i(x), f_i(x_0)] + D[f_i(x_0), f(x_0)] \\ &< \epsilon/3 + \epsilon/3 + \epsilon/3 = \epsilon. \end{aligned}$$

We are going to use the function space about to be introduced only in the case where Y is the metric space of real numbers. However, the generalization requires no extra work.

For any metric spaces we introduce $C(X, Y)$, the space of bounded continuous functions from X to Y . (We recall that $f: X \rightarrow Y$ is said to be bounded if $f(X)$ has finite diameter.) The metric is given by

$$D(f, g) = \sup_{x \in X} D[f(x), g(x)].$$

Observe that convergence in the metric of $C(X, Y)$ is precisely uniform convergence. Compare Examples 9(a), 11, and 12 in Appendix 1.

THEOREM 55. *Let X be a metric space, Y a complete metric space. Then $C(X, Y)$ is a complete metric space.*

Proof: Let $\{f_i\}$ be a Cauchy sequence in $C(X, Y)$. We get a limit for $\{f_i\}$ by a method typical of completeness proofs in function spaces: The desired limit function is located pointwise, and then the fact that $\{f_i\}$ is "uniformly Cauchy" is used to conclude the argument.

For each $x \in X$, $\{f_i(x)\}$ is plainly a Cauchy sequence in Y . Since Y is complete there is a limit, which we take as $f(x)$. We now verify that f_i converges uniformly to f . Given $\epsilon > 0$ we take N so that $D(f_i, f_j) < \epsilon/2$ for $i, j \geq N$. (As in the proof of Theorem 52, $\epsilon/2$ is used in this proof merely to take care of an impending degeneration of $<$ to \leq .) Fix x and i for the moment. Since we have $D[f_i(x), f_j(x)] < \epsilon/2$ for all j , we can pass to the limit by Theorem 37, and we get $D[f_i(x), f(x)] \leq \epsilon/2 < \epsilon$. Hence the convergence of f_i to f is uniform. We deduce from Theorem 54 that f is continuous. The final point needed is that f is bounded. This is easy: As soon as we have $D[f_i(x), f(x)] < \epsilon$ for all x we see that the diameter of $f(X)$ is at most the diameter of $f_i(X)$ increased by 2ϵ . Thus

$f \in C(X, Y)$, f_i converges to f in the metric of $C(X, Y)$, and the theorem is proved.

We can embed any metric space X in $C(X, R)$, the space of bounded real continuous functions on X .

THEOREM 56. *Let X be a metric space, and fix a point a in X . Assign to every $u \in X$ the real-valued function $f_u: X \rightarrow R$ given by $f_u(x) = D(u, x) - D(a, x)$. Then the map $u \rightarrow f_u$ is an isometry of X into $C(X, R)$.*

Proof: f_u is continuous by Exercise 2 in Section 4.4 and the continuity of the subtraction of real numbers. By Theorem 26, $|f_u(x)| \leq D(u, a)$, showing that f_u is bounded. Now let u and v be given; we have to prove that $D(u, v) = D(f_u, f_v)$. We have

$$(18) \quad \begin{aligned} D(f_u, f_v) &= \sup_{x \in X} |f_u(x) - f_v(x)| \\ &= \sup_{x \in X} |D(u, x) - D(v, x)|. \end{aligned}$$

Since

$$|D(u, x) - D(v, x)| \leq D(u, v)$$

by Theorem 26, we deduce $D(f_u, f_v) \leq D(u, v)$ from (18). Since $D(u, x) - D(v, x)$ takes the value $D(u, v)$ at the point v , we get the reverse inequality $D(u, v) \leq D(f_u, f_v)$. Hence the map $u \rightarrow f_u$ is an isometry, as required.

We are ready for the final touch on the question of completion. Let X be an arbitrary metric space. By Theorem 56, there is an isometry between X and a subset of $C(X, R)$. We can identify X with this subset. By Theorem 55 (and Theorem 47), $C(X, R)$ is complete. Now it is easy to complete a metric space once it has been placed inside a complete one: just take the closure (Theorem 48). We summarize:

THEOREM 57. *Any metric space has a completion.*

Whenever an example of a metric space comes up in a natural way, the question as to whether it is complete is almost certain to be of interest. More exactly, if it is complete, that will be an interesting theorem (or at least a worthwhile exercise); if it is incomplete, it is valuable to study the completion. Thus, the completion of Example 10 in Appendix 1 is called the p -adic numbers and plays a significant role in modern algebra and number theory. Example 9(b) is incomplete; forming its completion is (at least from one point of view) the essence of the theory of Lebesgue

integration. Example 9(c) is also incomplete. Its completion turns out to be isometric to Example 8 (Hilbert space).

EXERCISES

1. Prove that the space m (Example 7 in Appendix 1) is complete.
2. Let X be a dense subset of a metric space Y . Suppose that every Cauchy sequence in X converges to a point in Y . Prove that Y is complete. (This exercise is useful in the direct method of completing a metric space.)
3. Prove that the completion of an ultrametric space is ultrametric (see Exercise 6 in Section 4.1 for the definition).
4. Let X and Y be metric spaces. Assume that X is complete and that there exists $f: X \rightarrow Y$ which is one-to-one, onto, and continuous. Suppose that f^{-1} is uniformly continuous. Prove that Y is complete.
5. If every countable closed subset of a metric space M is complete, prove that M is complete.
6. If every closed ball of a metric space M is complete, prove that M is complete.
7. Let A and B be subsets of a metric space. If A and B are complete, prove that $A \cup B$ and $A \cap B$ are complete.
8. If a is a fixed point in a metric space M , prove that $D(a, x)$ is a uniformly continuous function from M to the real numbers.
9. If A is a fixed subset of a metric space M , prove that $D(A, x)$ is a uniformly continuous function from M to the real numbers.
10. Prove that the following three statements are equivalent for a metric space M :
 - (a) Every Cauchy sequence in M is ultimately constant.
 - (b) M is complete and discrete.
 - (c) Every subset is complete.
11. In the metric space of real numbers, give an example of a descending sequence of nonempty closed sets with empty intersection.
12. The following functions are continuous from the real numbers to the real numbers. Which are uniformly continuous?
 - (a) x^2 .
 - (b) $|x|$.
 - (c) $\frac{1}{1 + x^2}$.
13. Let R be the metric space of real numbers, and let $f: R \rightarrow R$ be a function which has a bounded derivative. Prove that f is uniformly continuous.

14. If, in a metric space, a sequence $\{x_i\}$ does not converge to y , prove that there exists a neighborhood U of y and a subsequence of $\{x_i\}$ lying outside U .
15. Let $\{x_i\}, y$ be a sequence and a point in a metric space. If every subsequence of $\{x_i\}$ has a subsequence converging to y , prove that $\{x_i\}$ converges to y .
- 16.* Prove that any open subset of a complete metric space is homeomorphic to a complete metric space. (*Hint*: If A is the complement of the open set, use the metric

$$D_1(x, y) = D(x, y) + \left| \frac{1}{D(x, A)} - \frac{1}{D(y, A)} \right|$$

- 17.* A metric space is said to be *convex* if for any $x \neq y$ there exists a point z different from x and y with $D(x, y) = D(x, z) + D(z, y)$. Prove that in a complete convex metric space there exists, for any x and y , a point t with $2D(x, t) = 2D(t, y) = D(x, y)$. (*Hint*: By transfinite induction steadily refine the point z to bring it closer to a midpoint, using the completeness at limit ordinals.)
18. Prove that there exist complete metric spaces of any cardinal number.

5.2 SEPARABILITY

We imitate for general metric spaces the way the rational numbers lie inside the real numbers as a countable dense subset.

DEFINITION. A metric space is *separable* if it possesses a countable dense subset.

As just noted, the space of real numbers is separable. So is Euclidean n -space: again we can take the points with rational coordinates. (Or we may take the point of view of product spaces, as in Section 6.1.) Not every metric space is separable. Indeed, in a discrete metric space M a dense subset has to be all of M . So to get an example which is not separable we need only take an uncountable discrete space, e.g., Example 4 in Section 4.1 for the case of an uncountable number of points.

It is illuminating to connect separability with two other properties.

DEFINITION. A collection $\{U_i\}$ of open subsets of a metric space M is called an *open base* if every open set in M is expressible as a union of U_i 's.

As an example of an open base we can take all the open sets in M . To be slightly more economical, we can take all open balls. In order to be still more economical the following theorem is helpful.

THEOREM 58. *Let $\{x_i\}$ be a dense subset of a metric space M . Consider the collection \mathfrak{U} of all open balls with rational radius and center one of the x_i 's. Then \mathfrak{U} is an open base for M .*

Proof: Let V be an open set in M . We must express V as a union of the prescribed open balls. We shall take $y \in V$ and prove that one of the members of \mathfrak{U} contains y and is contained in V . The first step is to pick an open ball with center y lying in V (Figure 12). Say the radius is r . We can suppose that r is rational. There exist elements x_i as close as we like to y . Pick one, say x_j , with $D(x_j, y) < r/2$. Then $S_{r/2}(x_j)$ contains y , and it lies inside V by the triangle inequality. Since y was an arbitrary element of V , this shows that V is a union of members of \mathfrak{U} .

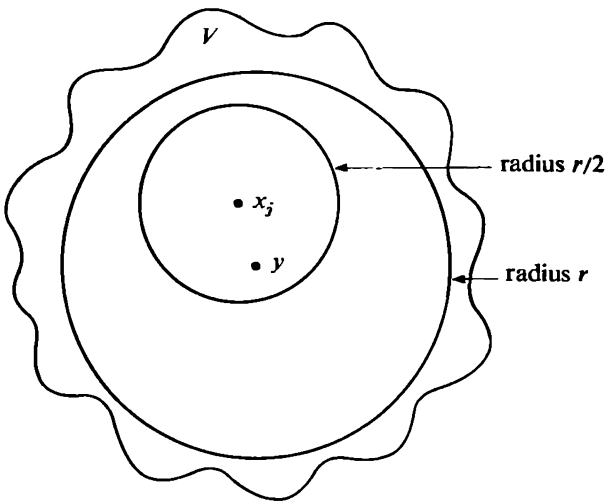


Figure 12

We next discuss coverings of a space. By an *open covering* of a metric space M we mean a collection of open sets whose union is M . A subcovering is obtained by deleting some of the sets in such a way that the surviving ones still cover M .

THEOREM 59. *The following three statements are equivalent in a metric space M :*

- (a) M is separable.
- (b) M has a countable open base.
- (c) Any open covering of M admits a countable subcovering.

Proof: (a) \Rightarrow (b): This implication is immediate from Theorem 58.

(b) \Rightarrow (c): Let $\mathfrak{u} = \{U_i\}$ be a countable open base for M , and let $\{V_j\}$ be an open covering of M . For each U_i that is contained in some V_j make (arbitrarily) a definite choice of such a V_j . Keep just these V_j 's, throwing all the rest away. Clearly we have kept only countably many V_j 's. We claim that these surviving V_j 's still cover M . To see this, we argue one point at a time. Pick any $x \in M$. x lies in some V_k , and that V_k in turn is a union of U_i 's. x must lie in at least one of these U_i 's, say U_m . Since $U_m \subset V_k$, some one of the V_j 's containing U_m (not necessarily V_k) was kept among the survivors, and it contains x .

(c) \Rightarrow (a): For each positive integer n , cover M by all open balls of radius $1/n$. Cut this down to a countable subcovering, and then pick a point at random in each of these surviving open balls. We obtain in this way a countable set, say T_n , with the following property: Any point of M lies within a distance $1/n$ of some point of T_n . The union of all the sets T_n ($n = 1, 2, 3, \dots$) is a countable dense subset.

A good illustration of the usefulness of having alternative criteria such as (b) and (c) is provided by Exercise 1. Given a separable metric space M and a subset A , it is not clear offhand how to relate to A the countable dense subset of M . However, property (b) or property (c) can be used in a straightforward way.

We turn to some cardinal number questions.

THEOREM 60. *A separable metric space has cardinal number at most c (the cardinal number of the set of real numbers).*

Proof: Let $\{x_i\}$ be a countable dense subset of the metric space M . Then (Theorem 36) every point of M is the limit of a convergent sequence of x_i 's. There are at most $\aleph_0^{\aleph_0}$ sequences in all (convergent or not), and $\aleph_0^{\aleph_0} = c$.

DEFINITION. A point x in a metric space M is a *condensation* point of M if every neighborhood of x contains an uncountable number of points of M .

THEOREM 61. *In an uncountable separable metric space M all but a countable number of points are condensation points.*

Proof: We first prove that there is at least one condensation point. Suppose the contrary. Then every point x in M has a countable open neighborhood N_x . The N_x 's are an open cover of M . By the criterion (c)

of Theorem 59, there is a countable subcovering. This gives us the contradiction that all of M is countable.

Next let Z be the set of all condensation points in M . We claim that $M - Z$ is countable. If not, we observe that $M - Z$ has a condensation point. This means, in the first instance, a condensation point relative to $M - Z$, but such a point is *a fortiori* a condensation point of M . This contradiction completes the proof.

THEOREM 62. *Any uncountable complete separable metric space has cardinal number c .*

Proof: That the cardinal number is at most c is asserted by Theorem 60. Conversely, the construction of c points in the given metric space M will be accomplished by repeated use of the following observation: M contains two disjoint uncountable closed balls of arbitrarily small radius. To see this, just take two condensation points (Theorem 61) and surround them by suitably small closed balls. Now start the construction with A_1 and A_2 disjoint uncountable closed balls of radius ≤ 1 . A_1 and A_2 are again separable and complete (Exercise 1 and Theorem 48). So the process can be repeated. In A_i ($i = 1, 2$) we can find disjoint uncountable closed balls A_{i1}, A_{i2} with radius $\leq \frac{1}{2}$. Repeated indefinitely, this procedure yields the following: For any sequence of 1's and 2's we have a descending sequence of closed balls with radii approaching 0. By Theorem 50 each such sequence clamps down on a point. We thus get c distinct points in M .

EXERCISES

1. Prove that any subspace of a separable metric space is separable. (*Hint:* In addition to Theorem 59, use Exercise 9 in Section 4.3.)
2. Let A_1, A_2, A_3, \dots be subsets of a metric space. Suppose that each A_i is separable. Prove that $\cup A_i$ is separable.
3. (a) Let X be a separable metric space. Let Y be a continuous image of X , i.e., there exists a continuous function from X onto Y . Prove that Y is separable.
(b) Is part (a) true if "separable" is replaced by "complete"?
- 4.* Let M be a metric space. Prove that M is separable if and only if every collection of disjoint open sets of M is countable.
- 5.* Let M be a separable metric space. Let \mathcal{C} be a collection of closed subsets of M such that whenever $A_1 \supset A_2 \supset A_3 \supset \dots$ is a descending sequence of elements of \mathcal{C} then the intersection $\cap A_i$ lies in \mathcal{C} . Prove that \mathcal{C} has a minimal element. (*Hint:* Use Zorn's lemma and Theorem 59.)

- 6.* Let d be an infinite cardinal. Prove that the following statements are equivalent for a metric space M :
- M has a dense subset of cardinal $\leq d$.
 - M has an open base of cardinal $\leq d$.
 - Any covering of M by open subsets admits a subcovering of cardinal $\leq d$.
7. Prove that in any metric space the closure of a separable subspace is separable.
8. Prove that any uncountable complete metric space has at least the cardinal number c .
- 9.* Prove that the metric space of all bounded sequences of real numbers (Example 7 in Appendix 1) is not separable.
10. Prove that the following is equivalent to separability of a metric space M : for any collection $\{F_i\}$ of closed sets in M , if every countable intersection of F_i 's is nonempty, then $\bigcap F_i$ is nonempty.

5.3 COMPACTNESS

In this section we shall treat the most important concept to be introduced in our treatment of metric spaces: compactness. In Theorem 64 we shall note the equivalence of two "sequential" versions of compactness, and in Theorem 72 we pass to the characterization of compactness by open coverings.

We first need to note a connection between limit points and convergent subsequences. Theorem 63 is a slight variant of results like Theorem 36; it could also be deduced from Exercise 18 in Section 4.3.

THEOREM 63. *Let $\{x_i\}$ be a sequence of distinct points in a metric space M . Write A for the (unordered) set of x_i 's. Then for any point x lying in the closure of A but not in A there exists a subsequence of $\{x_i\}$ converging to x .*

Proof: Any neighborhood of x contains an infinite number of elements of A . We pick the desired subsequence of $\{x_i\}$ by successive steps; in the n th step we pick an element of $S_{1/n}(x)$ lying beyond the already selected elements in the sequence $\{x_i\}$.

THEOREM 64. *The following statements are equivalent in a metric space:*

- Every sequence has a convergent subsequence.
- Every infinite set has a limit point. (See the end of Section 4.3 for the definition of a limit point.)

Proof: That (a) implies (b) is obvious. In proving that (b) implies (a) we can assume that the given sequence $\{x_i\}$ consists of different elements. Then $\{x_i\}$ has, by hypothesis, a limit point x . We can suppose that x is not one of the x_i 's (if $x_j = x$ just delete x_j). The existence of a convergent subsequence now follows from Theorem 63.

DEFINITION. A metric space is *compact* if it satisfies either (hence both) of the conditions in Theorem 64.

THEOREM 65. *Any closed interval on the real line (i.e. the set of all x satisfying $a \leq x \leq b$) is compact.*

Proof: Only a small modification of the proof of Theorem 47 is needed, and we leave it to the reader.

The next theorem is the analogue of Theorem 48.

THEOREM 66. *A closed subset of a compact metric space is compact.*

Proof: Let M be a compact metric space and A a closed subset of M . Let $\{a_i\}$ be a sequence of elements in A . There is a subsequence converging to an element x of M . Since A is closed, $x \in A$.

Compactness is stronger than completeness.

THEOREM 67. *Any compact metric space M is complete.*

Proof: Let $\{x_i\}$ be a Cauchy sequence in M . Since M is compact there is a subsequence converging, say, to y . By Theorem 44, the whole sequence converges to y .

Remark: As a corollary of Theorems 67 and 49 we get that any compact subset of a metric space is closed. This is, of course, easily seen directly as well.

THEOREM 68. *A compact metric space M has a finite diameter. The diameter is attained for two suitable points of M .*

Proof: Let t be the diameter of M . For the moment we allow $t = \infty$. We can find points $x_i, y_i \in M$ such that $D(x_i, y_i)$ approaches t . Since M is compact, we can drop down to convergent subsequences; so (after changing notation) we may assume that $x_i \rightarrow x, y_i \rightarrow y$. Then (Theorem 37) $D(x_i, y_i) \rightarrow D(x, y)$, so t is finite and equals $D(x, y)$.

As a consequence of these results we see that *a subset of the real line is compact if and only if it is closed and bounded*. This is true equally well in Euclidean space, but we postpone the discussion till the section on product spaces (Section 6.1).

It is easy to see that if a closed set of real numbers is bounded from above, then it contains its least upper bound. From this and Exercise 1 we deduce the following corollary:

THEOREM 69. *A continuous real function on a compact metric space is bounded and attains its upper and lower bounds.*

We need the fact that compactness implies separability.

THEOREM 70. *Any compact metric space M is separable.*

Proof: For each positive integer n , let T_n be a subset of M maximal with respect to the property that the distance between any two elements of T_n is at least $1/n$. (We can of course get such a maximal set by Zorn's lemma. But compactness makes it clear that no infinite set of this kind can exist—it would be impossible to extract a convergent sequence. So T_n is finite.) Note that any element of M lies within $1/n$ of some member of T_n . The union of the T_n 's is a countable dense set in M .

We now head toward the major characterization of compactness given in Theorem 72. We first prove Theorem 71, partly to break up the proof of Theorem 72, and partly to exhibit the way Theorem 71 is an analogue of Theorem 50.

THEOREM 71. *The following statements are equivalent for a metric space M :*

- (a) *M is compact.*
- (b) *If $F_1 \supset F_2 \supset F_3 \supset \dots$ is a descending sequence of nonempty closed sets in M , then $\bigcap F_i$ is nonempty.*

Remark: Observe that the assumption $\text{diam}(F_i) \rightarrow 0$, which occurred in Theorem 50, is not pertinent here.

Proof: (a) \Rightarrow (b): Pick any $x_i \in F_i$. A subsequence of $\{x_i\}$ converges, say to y . For arbitrarily large j we have that y is the limit of a sequence of elements in F_j , hence lies in F_j . Therefore, $y \in \bigcap F_i$.

(b) \Rightarrow (a): Suppose on the contrary that M is not compact. Then there exists a sequence $\{x_i\}$ with no convergent subsequence, and we can assume that it consists of distinct elements. For $i = 1, 2, 3, \dots$, let

$F_i = \{x_i, x_{i+1}, \dots\}$. It follows from Theorem 63 that F_i is closed. Furthermore, the F 's form a descending sequence of closed sets with empty intersection, and this contradicts our hypothesis.

THEOREM 72. *The following statements are equivalent for a metric space M :*

- (a) M is compact.
- (b) Any open covering of M admits a finite subcovering.
- (c) If a family of closed sets $\{F_i\}$ in M is such that any finite number of F_i 's have nonempty intersection, then $\bigcap F_i$ is nonempty.

Proof: We can pass back and forth between (b) and (c) by complementation (compare Exercise 10 in Section 5.2). Also (c) \Rightarrow (a) by Theorem 71. It remains to check (a) \Rightarrow (b).

(a) \Rightarrow (b): By Theorems 70 and 59, we can reduce the open covering to a countable open covering, say by sets U_1, U_2, \dots . Suppose that no finite number of these cover M . Write F_n for the complement of $U_1 \cup U_2 \cup \dots \cup U_n$. Then the F 's form a decreasing sequence of closed nonempty sets with empty intersection. This contradicts Theorem 71.

Uniform continuity is automatic on a compact space.

THEOREM 73. *Let f be a continuous function from a compact metric space X to a metric space Y . Then f is uniformly continuous.*

Proof: We make an indirect proof. Suppose there is an ϵ that cannot be matched by a δ , however small. By taking a sequence of δ 's approaching 0 we get sequences x_n, y_n with $D(x_n, y_n) \rightarrow 0$ and $D[f(x_n), f(y_n)] \geq \epsilon$. A subsequence of $\{x_n\}$ converges, say to x . After a change of notation, we may assume that the entire sequence $\{x_n\}$ converges to x . Since $D(x_n, y_n) \rightarrow 0$, $\{y_n\}$ also converges to x . From the continuity of f , we deduce that $f(x_n)$ and $f(y_n)$ both converge to $f(x)$. This contradicts the fact that $D[f(x_n), f(y_n)] \geq \epsilon$ for all n .

The condition that a metric space have a compact completion has several useful variants. They are presented in Theorem 74.

THEOREM 74. *The following conditions are equivalent for a metric space M :*

- (a) The completion M^* of M is compact.
- (b) Every sequence has a Cauchy subsequence.
- (c) For every $\epsilon > 0$ there exists a finite set T_ϵ such that every point in M is within ϵ of a member of T_ϵ .

- (d) For every $\epsilon > 0$, M is the union of a finite number of sets, each of diameter $< \epsilon$.
- (e) Any infinite subset of M contains an infinite subset of diameter $< \epsilon$.

Proof: We shall prove (a) \Leftrightarrow (b), and then run (b), (c), (d), and (e) in a circle.

(a) \Rightarrow (b): A sequence of elements of M has a subsequence converging to a point of M^* . By Theorem 43, this subsequence is Cauchy.

(b) \Rightarrow (a): Let $\{y_i\}$ be a sequence in M^* . We can, for each i , find x_i in M with $D(x_i, y_i) < 1/i$. From the fact that $\{x_i\}$ has a Cauchy subsequence it follows readily that $\{y_i\}$ does too. Since M^* is complete, this shows that M^* is compact.

(b) \Rightarrow (c): Take T_ϵ to be a subset of M maximal with respect to the property that any two members of T_ϵ are at distance $\geq \epsilon$. T_ϵ cannot be infinite, for such an infinite set can be exhibited as a sequence with no Cauchy subsequence.

(c) \Rightarrow (d): Given $\epsilon > 0$, let u_1, \dots, u_n be a finite set such that every point of M is within $\epsilon/2$ of one of u_1, \dots, u_n . Then M is the union of the open balls $S_{\epsilon/2}(u_1), \dots, S_{\epsilon/2}(u_n)$.

(d) \Rightarrow (e): Let A be an infinite subset of M and let ϵ be given. We have $M = B_1 \cup \dots \cup B_n$ where each B_i has diameter $< \epsilon$. Since

$$A = (A \cap B_1) \cup \dots \cup (A \cap B_n),$$

at least one of the sets $A \cap B_i$ must be infinite. Of course $A \cap B_i$ *a fortiori* has diameter $< \epsilon$.

(e) \Rightarrow (b): Let $\{x_i\}$ be a sequence in M . We prove that $\{x_i\}$ has a Cauchy subsequence by using Cantor's diagonal device. We can assume all the x_i 's to be distinct. By using our hypothesis (e) we can pick a subsequence such that all mutual distances are less than 1. We write this subsequence as

$$(19) \quad x_{11}, x_{12}, x_{13}, \dots$$

Now use (e) again to pick a subsequence of (19) with all mutual distances less than $\frac{1}{2}$. Our notation for this is

$$(20) \quad x_{21}, x_{22}, x_{23}, \dots$$

At the third step pick a subsequence of (20) with all mutual distances less than $\frac{1}{3}$. This procedure is continued indefinitely. The result is a doubly infinite array

$$\begin{array}{ccccccc} x_{11}, & x_{12}, & x_{13}, & \dots, & x_{1n}, & \dots & \\ x_{21}, & x_{22}, & x_{23}, & \dots, & x_{2n}, & \dots & \\ & & & \dots & & & \\ x_{n1}, & x_{n2}, & x_{n3}, & \dots, & x_{nn}, & \dots & \\ & & & \dots & & & \end{array}$$

The diagonal sequence

$$x_{11}, x_{22}, x_{33}, \dots, x_{nn}, \dots$$

is a Cauchy subsequence of the original sequence $\{x_i\}$.

The classical name for a space satisfying any (hence all) of the conditions in Theorem 74 is *totally bounded*. The designation *precompact* is being used by many authors, and the change is a good idea.

We conclude this section with a theorem on Lebesgue numbers of open coverings.

DEFINITION. Let $\{U_i\}$ be an open covering of a metric space M . A *Lebesgue number* for $\{U_i\}$ is a positive real number ϵ with the following property: If A is any subset of M with $\text{diam}(A) < \epsilon$, then A is contained in some U_i .

THEOREM 75. *Any open covering of a compact metric space has a Lebesgue number.*

Proof: Assume on the contrary that $\{U_i\}$ is an open covering of the compact metric space M , and that $\{U_i\}$ admits no Lebesgue number. It follows that there exists a sequence $\{A_n\}$ of sets such that $\text{diam}(A_n) \rightarrow 0$ but each A_n is not contained in any U_i . Pick $x_n \in A_n$. The sequence $\{x_n\}$ has a convergent subsequence; after a change of notation we may assume that $x_n \rightarrow x$. Suppose that $x \in U_i$ and that the open ball $S_{2r}(x)$ is contained in U_i . For sufficiently large n we have that $D(x_n, x) < r$ and $\text{diam}(A_n) < r$. But then for any y in A_n we have

$$D(x, y) \leq D(x, x_n) + D(x_n, y) < r + r = 2r,$$

showing that $A_n \subset S_{2r}(x) \subset U_i$, a contradiction.

Theorem 75 has useful applications in algebraic topology. It can also be used to give an alternative proof of Theorem 73 (see Exercise 15).

EXERCISES

1. Prove that a continuous image of a compact metric space is compact.
2. Prove that a continuous function from a compact metric space to any metric space is closed (i.e., it sends closed sets into closed sets).
3. Let f be a continuous one-to-one mapping of a compact metric space X onto a metric space Y . Prove that $f^{-1}: Y \rightarrow X$ is continuous. (*Hint:* Use the preceding exercise and Exercise 10(r) in Section 4.4.)

- 4.* (a) Let $\{f_i\}$ and f be continuous real functions on a compact metric space M . Prove that f_i converges uniformly to f if and only if the following is true: When a sequence x_1, x_2, x_3, \dots converges to x then $f_1(x_1), f_2(x_2), f_3(x_3), \dots$ converges to $f(x)$.
 (b) Give an example to show that the assumption of compactness cannot be omitted in the "only if" portion of part (a).
5. Let A and B be disjoint subsets of a metric space. Suppose that A is closed and B is compact. Prove that the distance between A and B is positive. (See Exercise 10 in Section 4.3 for the definition of the distance between two sets. Compare with part (b) of that exercise.)
6. A real-valued function f on a metric space M is called *upper semicontinuous* if, for each real number t , the set of all x in M with $f(x) \geq t$ is closed in M . Prove that if M is compact, an upper semicontinuous function on M is bounded above and attains its least upper bound.
7. Let $\{f_i\}$ be a sequence of isometries of a metric space M onto itself. Suppose that f_i converges to f pointwise, i.e., for every $x \in M$, $f_i(x)$ converges to $f(x)$. Prove that f is an isometry. If M is compact, prove that the convergence of f_i to f is uniform.
- 8.* Let f_1, f_2, f_3, \dots be continuous real functions on a compact metric space. Suppose that $f_1 \geq f_2 \geq f_3 \geq \dots$ and that the f_i 's converge pointwise to a continuous function f . Prove that the convergence is uniform.
9. Prove that in a metric space a subset is closed if and only if its intersection with every compact set is closed.
10. Let M be a compact metric space. Let A be a set of real functions on M . Assume that A is closed under multiplication and that for any $x \in M$ there exists a function in A vanishing in a neighborhood of x . Prove that A contains the zero function. (The multiplication meant is pointwise: $(fg)(x) = f(x)g(x)$.)
11. Let A and B be subsets of the real numbers. Write $A + B$ for the set of all $a + b$ ($a \in A, b \in B$). If A and B are closed, does it follow that $A + B$ is closed? Does this follow if in addition one of the two is bounded? \checkmark_{e}
- 12.* Let M be a compact metric space and $f: M \rightarrow M$ a function such that $D[f(x), f(y)] \geq D(x, y)$ for all x, y in M . Prove that f is an isometry of M onto itself. (Hint: For any $a \in M$, denote by a_n the result of applying f n times to a . For $a, b \in M$, extract "simultaneous" convergent subsequences of $\{a_n\}$ and $\{b_n\}$. In this way, get k with $D(a, a_k)$ and $D(b, b_k)$ arbitrarily small, showing that $D[f(a), f(b)] > D(a, b)$ is impossible. Thus f is an isometry. Furthermore, since $D(a, a_k)$ is arbitrarily small, the range of f is dense. Observe that the range of f is closed.)
- 13.* Let $f: M \rightarrow M$ be a mapping on a compact metric space which is one-to-one, onto, and satisfies $D[f(x), f(y)] \leq D(x, y)$ for all $x, y \in M$. Prove that f is an isometry. (Hint: Use Exercise 12.)
14. Let M be the metric space consisting of $1, \frac{1}{2}, \frac{1}{3}, \dots, \frac{1}{n}, \dots$ (not including 0). Exhibit an open covering of M with no Lebesgue number.

15. Use Theorem 75 to give an alternate proof of Theorem 73. (*Hint:* Given ϵ , and a point $x \in X$, pick the γ that corresponds to $\epsilon/2$ in the continuity of f at x . Cover X by the open balls $S_\gamma(x)$, and let δ be a Lebesgue number for this covering.)
- 16.* For a metric space M , prove the following statements equivalent:
 - (a) M is compact.
 - (b) Every infinite open covering of M admits a proper subcovering (i.e., in an open covering of M by an infinite number of open sets at least one of the sets can be discarded with the remaining sets still covering M).
17. If A and B are compact subsets of a metric space, prove that $A \cup B$ is compact.

6

Additional Topics

6.1 PRODUCT SPACES

Let X and Y be metric spaces. In this discussion it will be advisable to use different letters (say D_1 and D_2) for the metrics on X and Y . We now seek to put a metric on the Cartesian product $X \times Y$. The difficulty in the discussion is an embarrassment of riches. There are too many plausible ways of putting a metric on $X \times Y$. At the very least there are the three choices suggested by Example 3 in Section 4.1: $\sqrt{D_1^2 + D_2^2}$, $\text{Max}(D_1, D_2)$, and $D_1 + D_2$.

It is a fact that for almost all purposes the choice is not important. To see this in a systematic way, one should note that the identity map on $X \times Y$, from one of the above metrics to another, is uniformly continuous in both directions; then observe that the major properties of metric spaces are invariant under such a map. (*Remark:* One exception is the diameter.)

Rather than pursue the full details, let us take $\text{Max}(D_1, D_2)$ as our "official" choice for a metric on $X \times Y$. In detail: If $u_1 = (x_1, y_1)$, $u_2 = (x_2, y_2)$ are two points of $X \times Y$, we set

$$D(u_1, u_2) = \max [D_1(x_1, x_2), D_2(y_1, y_2)].$$

It should be noted that this does not give us the usual metric in the Euclidean plane.

We make a number of observations, omitting the easy proofs.

1. If U is open in X and V is open in Y , then $U \times V$ is open in $X \times Y$. However, not every open set of $X \times Y$ is obtainable this way. Exactly the same remark can be made for closed sets.

2. Convergence in $X \times Y$ means pointwise convergence: $(x_i, y_i) \rightarrow (x, y)$ if and only if $x_i \rightarrow x$ and $y_i \rightarrow y$.

3. Let Z be a third metric space. A continuous function from Z to $X \times Y$ is simply a combination of a continuous function from Z to X with one from Z to Y . However, a continuous function from $X \times Y$ to Z is a function of two variables which is jointly continuous, rather than merely separately continuous in its variables.

4. Now we are able to make an overdue remark: The distance function on a metric space M is a continuous function (from $M \times M$ to R).

5. The Cartesian product of a finite number of metric spaces can be defined by iterating the product of two spaces, or in one fell swoop. It makes no real difference which is done. All the statements in Theorem 76 are valid for the product of any finite number of metric spaces.

6. On the other hand, the Cartesian product of an infinite number of metric spaces presents some difficulties. It will be better to delay the discussion to Appendix 3, where the matter will be surveyed in the general context of topological spaces.

We assemble in a single theorem a number of statements, each asserting that a product space $X \times Y$ has a certain property if and only if both X and Y do.

THEOREM 76. *Let X and Y be metric spaces. Each of the following properties holds for $X \times Y$ if and only if it holds for both X and Y : completeness, compactness, separability, precompactness, finiteness of the diameter.*

The proof is quite routine and is left to the reader.

EXERCISES

1. Let X and Y be any metric spaces. Prove that the natural projection of $X \times Y$ onto X is open (i.e., it sends open sets into open sets).
2. Prove that a metric space A is compact if and only if for each metric space B the natural projection of $A \times B$ onto B carries closed sets into closed sets.

3. (a) Let f be a continuous map from a metric space X to a metric space Y . Prove that the graph of f is a closed subset of $X \times Y$.
 (b) If X is compact and the graph of f is closed, prove conversely that f is continuous.
 (c) Is compactness of X needed in part (b)?
4. Prove: In the Cartesian product of a metric space X with itself, the diagonal (the set of all (x, x) with $x \in X$) is closed.

6.2 A FIXED-POINT THEOREM AND AN APPLICATION

A fixed-point theorem is an assertion of the following kind: For a suitable mapping f on a suitable set X there exists a point $x \in X$ with $f(x) = x$. For metric spaces there is a very simple fixed-point theorem for mappings which might appropriately be called *strict contractions*.

THEOREM 77. *Let M be a complete metric space and let $f: M \rightarrow M$ be a mapping with the following property: There exists a real number $k < 1$ such that $D[f(x), f(y)] \leq kD(x, y)$ for all $x, y \in M$. Then f has a unique fixed point.*

Proof: The uniqueness is immediate (and does not require the completeness of M); for if $f(x) = x$ and $f(y) = y$ with $x \neq y$, then

$$D[f(x), f(y)] = D(x, y),$$

contrary to the hypothesis.

To prove the existence of a fixed point, we construct a Cauchy sequence by iterating f . Start with any point z in M and write $z_1 = z$, $z_2 = f(z_1)$, \dots , $z_n = f(z_{n-1})$, etc. We claim that $\{z_i\}$ is a Cauchy sequence. We have

$$(21) \quad D(z_{n+1}, z_n) = D[f(z_n), f(z_{n-1})] \leq kD(z_n, z_{n-1}).$$

By iterating (21) and using the triangle inequality, we have, for $m > n$,

$$\begin{aligned} D(z_m, z_n) &\leq D(z_m, z_{m-1}) + D(z_{m-1}, z_{m-2}) + \dots + D(z_{n+1}, z_n) \\ &\leq (1 + k + k^2 + \dots + k^{m-n-1})D(z_{n+1}, z_n) \\ &= \frac{1 - k^{m-n}}{1 - k} D(z_{n+1}, z_n) \\ &\leq \frac{k^{n-1}(1 - k^{m-n})}{1 - k} D(z_2, z_1). \end{aligned}$$

Now by taking n sufficiently large we can arrange for

$$\frac{k^{n-1}(1 - k^{m-n})}{1 - k} < \epsilon$$

to be as small as we like. We have sustained the claim that $\{z_i\}$ is a Cauchy sequence. (*Remark:* The last portion of the argument can be replaced by a reference to the convergence of a geometric series with ratio less than 1.)

Since M is complete, the sequence $\{z_i\}$ converges to a limit w , and since f is continuous, $f(z_i) \rightarrow f(w)$. But $f(z_i) = z_{i+1}$. Hence $f(w) = w$, and w is the desired fixed point.

We shall sketch a way of deducing from Theorem 77 an existence theorem for solutions of differential equations.

Let $f(x, y)$ be a continuous real-valued function of two real variables, defined in a neighborhood of the origin $(0, 0)$. We discuss solutions of the differential equation $y' = f(x, y)$ which satisfy the boundary condition $y(0) = 0$. Our discussion is "local." By this we mean that existence and uniqueness are to be proved only in a suitably small neighborhood of $x = 0$, and there is no investigation as to how far out the solution goes. Subject to a suitable Lipschitz condition on f , we shall reduce the problem to Theorem 77, and thereby obtain existence and uniqueness of a solution.

The first step is to replace the differential equation

$$(22) \quad y' = f(x, y)$$

by the integral equation

$$(23) \quad y(x) = \int_0^x f[t, y(t)] dt.$$

It is simple to pass back and forth between (22) and (23) by integration and differentiation.

Let f be defined and continuous on the rectangle R given by $|x| \leq A$ and $|y| \leq B$. Suppose that $|f| \leq K$ on R and that

$$(24) \quad |f(x, y_1) - f(x, y_2)| \leq L|y_1 - y_2|$$

holds on R . Here K and L are fixed positive real numbers. The inequality (24) is described by saying that, on R , f satisfies a Lipschitz condition on y .

By decreasing A , if necessary, we can arrange that $AK \leq B$ and $AK < 1$. We shall suppose that this has been done.

Let M be the metric space of all continuous real functions ϕ defined on the closed interval $[-A, A]$ and satisfying $|\phi| \leq B$. We make M into a metric space by using the distance function derived from the sup norm:

$$\|\phi\| = \sup_{|x| \leq A} |\phi(x)|.$$

On M we define a mapping $\phi \rightarrow \psi$ by sending ϕ into the function ψ given by

$$(25) \quad \psi(x) = \int_0^x f[t, \phi(t)] dt.$$

First we have to see that this makes sense. For $|t| \leq A$ we have $|\phi(t)| \leq B$. Hence the integrand in (25) is defined for $|x| \leq A$. The continuity of ϕ and f tells us that the integral exists. It is easy to see that ψ is continuous. We have the estimate

$$|\psi(x)| \leq |x| |f| \leq AK \leq B.$$

Hence $\psi \in M$. Next, if ϕ and ϕ_1 in M are sent into ψ and ψ_1 we have

$$(26) \quad |\psi(x) - \psi_1(x)| \leq \int_0^x |f[t, \phi(t)] - f[t, \phi_1(t)]| dt.$$

By (24), the integrand in (26) is dominated by $K|\phi(t) - \phi_1(t)|$. Hence

$$D(\psi, \psi_1) \leq AKD(\phi, \phi_1).$$

Since $AK < 1$ we have that the map $\phi \rightarrow \psi$ is a strict contraction on the metric space M . Since M is complete, we have a unique fixed point for the map by Theorem 77. This says that (23) has a unique solution for $|x| \leq A$.

The fixed-point technique can also be used to prove the implicit function theorem of multivariable calculus, in a coordinate-free form that applies to infinite-dimensional spaces. (See S. Lang's *Introduction to Differentiable Manifolds*, pages 12–13.)

EXERCISES

- Let f and g be commuting maps on a set. Prove that the set of fixed points of g is invariant under f .
 - Let f be a map on a set. Suppose that some power f^k has a unique fixed point x . Prove that x is fixed under f , and is the only fixed point of f .
 - Let f be a map (not necessarily continuous) on a complete metric space M . Suppose that some power of f is a strict contraction. Prove that f has a unique fixed point.
- Let f and g be commuting maps on a complete metric space. Suppose that f is a strict contraction (g need not be continuous). Prove that there exists a unique joint fixed point for f and g .
- * Let M be a compact metric space and let $f: M \rightarrow M$ be a function such that $D[f(x), f(y)] < D(x, y)$ for all $x, y \in M$, $x \neq y$. Prove that f has a unique fixed point. (*Hint*: Take z_i as in the proof of Theorem 77. A subsequence z_{n_i} converges to w . Suppose that $f(w) \neq w$. There are neighborhoods U, V of $w, f(w)$ and a real number $k < 1$ such that $D[f(x), f(y)] \leq kD(x, y)$ for $x \in U, y \in V$. For large i and $j > i$,

$$D[z_{n_j}, f(z_{n_i})] \leq k^{i-j} D[z_{n_i}, f(z_{n_i})].$$

So the left side is arbitrarily small for large j . This is a contradiction, since z_{n_j} and $f(z_{n_i})$ come close to w and $f(w)$. This proof is due to Edelstein.

See Belluce and Kirk, *Can. Math. Bull.* **12** (1969), 481–491, for recent results of this kind and references to earlier literature.)

- 4.* Let M be a complete metric space. Let f_i ($i = 1, 2, 3, \dots$) and f be strict contractions on M , and suppose that f_i converges uniformly to f . Let x_i, x be the fixed points of f_i, f . Prove that $x_i \rightarrow x$.

6.3 CATEGORY

The concepts of first and second category are ways of describing, in a certain sense, the size of a metric space. They are based in turn on the concept of a nowhere dense set.

DEFINITION. A subset A of a metric space M is *nowhere dense* in M if the closure of A does not contain a nonempty open set.

Remark: It is evidently equivalent to say that the closure of A does not contain a ball (open or closed).

We offer some examples of sets which are nowhere dense and others which are not. First consider a subset A consisting of just one point x . Evidently A is nowhere dense if and only if x is not isolated. Let us bypass this distorting influence of isolated points by assuming that our metric space M has no isolated points. Then of course any finite set is nowhere dense. A countable set may or may not be nowhere dense. For instance, let M be the closed unit interval $[0, 1]$. The subset A of rational numbers is a countable subset which is far from being nowhere dense, indeed, the closure of A is all of M . On the other hand, if A consists of a convergent sequence (with or without its limit point), then A is nowhere dense in $[0, 1]$.

It is a fact that uncountable nowhere dense sets also exist in $M = [0, 1]$, but it is not easy to exhibit an example. Perhaps the simplest is the “middle third” set discovered by Cantor. Delete from M the open middle third interval $(\frac{1}{3}, \frac{2}{3})$. From each of the remaining closed intervals $[0, \frac{1}{3}]$ and $[\frac{2}{3}, 1]$ delete the open middle thirds $(\frac{1}{9}, \frac{2}{9})$ and $(\frac{7}{9}, \frac{8}{9})$. The process is continued forever, and we denote by A the subset of M that survives. Then it is not difficult to see that A is closed, that it has the cardinal number of the continuum, and that it is nowhere dense.

DEFINITION. A subset of a metric space is of the *first category* if it is expressible as a countable union of nowhere dense sets; otherwise it is of the *second category*.

Theorem 79 is the main result concerning category. We first note Theorem 78.

THEOREM 78. *Let A be a nowhere dense set in a metric space M and let U be a nonempty open set in M . Then A is disjoint from some ball contained in U .*

Proof: Assume the contrary. Let x be a point in U . Any sufficiently small ball with center x is contained in U and therefore intersects A . Hence x lies in the closure \bar{A} of A . Thus \bar{A} contains U , a contradiction of the hypothesis that A is nowhere dense.

THEOREM 79. *A complete metric space is of the second category.*

Proof: Suppose on the contrary that the complete metric space M is a countable union $\cup A_i$ of nowhere dense sets. We begin a construction which leads us to a contradiction. By Theorem 78, A_1 is disjoint from a ball; we can take it to be a closed ball S_1 of radius ≤ 1 . Then we apply Theorem 78 to A_2 and an open ball contained in S_1 . We can thus get a closed ball S_2 with $S_2 \subset S_1$, $S_2 \cap A_2$ empty, and radius $(S_2) \leq \frac{1}{2}$. Continuing in this way, we get a descending sequence $\{S_i\}$ of closed balls with $S_n \cap A_n$ empty, and radius $(S_n) \leq 1/n$. By Theorem 50 there is a point $y \in \cap S_i$. But y lies in none of the A_i 's, contradicting the hypothesis $M = \cup A_i$.

The Polish mathematician Mazur invented a game to illustrate the idea of category. In this game two players begin by dividing the unit interval between them in some fashion. The first player picks a closed interval of length $\leq \frac{1}{2}$; the second then selects a closed interval of length $\leq \frac{1}{3}$ inside the already chosen interval; the first player now picks a closed interval of length $\leq \frac{1}{4}$ inside the preceding one, etc. These intervals clamp down on a unique point. The player having this point in his set wins.

Informal exercise: If one of the players is so unfortunate as to have a set of the first category, prove that his opponent can force a win.

EXERCISES

1. Prove that the union of a finite number of nowhere dense sets is nowhere dense.
2. Prove that any subset of a set of the first category is again of the first category, and that any superset of a set of the second category is again of the second category.

3. Prove that a nowhere dense set cannot contain any isolated points.
4. Let M be a countable metric space. Prove that M is of the first category if and only if it contains no isolated points.
5. Prove that a countable infinite complete metric space has an infinite number of isolated points.
6. Let M be a complete metric space and let $\{A_i\}$ be a countable collection of dense open subsets of M . Prove that $\bigcap A_i$ is dense.
7. Let X and Y be metric spaces, let A be a nowhere dense subset of X , and let B be a nowhere dense subset of Y . Prove that $A \times B$ is nowhere dense in $X \times Y$.
8. If X and Y are metric spaces of the first category, prove that $X \times Y$ is of the first category.
- 9.* In the metric space of real numbers, prove that the subset of irrational numbers is not expressible as a union of a countable number of closed sets.

Appendixes

1

Examples of Metric Spaces

The discussion of examples will be continued where it left off in Section 4.1.

Above all we wish to extend metric spaces like those in Example 3 of Section 4.1 to higher finite dimensions and to the infinite-dimensional case. The treatment will be facilitated by introducing the concept of a normed abelian group.

DEFINITION. A *normed abelian group* is an abelian group G together with a real-valued function on G , written $\|x\|$, which satisfies:

- I.' $\|0\| = 0$.
- II.' $\|x\| > 0$ for $x \neq 0$.
- III.' $\|x\| = \|-x\|$.
- IV.' $\|x + y\| \leq \|x\| + \|y\|$.

IV', like IV, is called the *triangle inequality*.

Remarks: 1. Some information on abelian groups can be found in Appendix 2.

2. By further specialization one passes from normed abelian groups to normed linear spaces, in which a real vector space structure is postulated,

satisfying in addition $\|cx\| = |c| \|x\|$ for any real number c . Normed linear spaces furnish the basic framework for the branch of mathematics called *functional analysis*. For our purposes it is a little simpler to work with normed abelian groups; furthermore, one of the examples (number 10, which follows) is not a normed linear space.

The following theorem shows how to get a metric space from any normed abelian group.

THEOREM. *Let G be a normed abelian group, relative to the norm $\|x\|$, $x \in G$. Define $D(a, b)$ to be $\|a - b\|$. Then G is a metric space relative to D .*

Proof: In the transition from I'-IV' to I-IV (the metric space postulates as given in Section 4.1), only IV (the triangle inequality) needs attention. Since $a - c = a - b + b - c$, we have, by IV',

$$\|a - c\| \leq \|a - b\| + \|b - c\|,$$

which gives us IV.

In each of the examples that follow, an abelian group is given and a proposed norm for it. Only the axiom IV' deserves any discussion, and there will be occasional remarks concerning its verification. We continue the numbering of examples from Section 4.1.

Example 5: This trio of examples extends those in Example 3 to any finite dimension. Let R_n be Euclidean n -space, i.e. the set of all ordered n ples of real numbers. For a point $a = (a_1, \dots, a_n)$ of R_n , the three norms are

- (a) $\|a\| = \sqrt{a_1^2 + \dots + a_n^2}$,
- (b) $\|a\| = \max(|a_1|, \dots, |a_n|)$,
- (c) $\|a\| = |a_1| + \dots + |a_n|$.

Example 6: Let l_1 denote the set of all sequences of real numbers $a = \{a_i\}$ with $\sum_{i=1}^{\infty} |a_i|$ finite (i.e. convergent). We take this sum to be the norm $\|a\|$. This is the infinite-dimensional analogue of Examples 3(c) and 5(c). There is one point to be carefully noted at once: We do not yet know that l_1 is an abelian group, that is, the closure of l_1 under addition is in doubt. As a matter of fact, we can prove this simultaneously with the triangle inequality. First we work with finite sums. Let $a = \{a_i\}$ and $b = \{b_i\}$ be elements of l_1 , and write $c_i = a_i + b_i$. By iterated use of the ordinary triangle inequality for real numbers, we get

$$(27) \quad |c_1| + \dots + |c_n| \leq |a_1| + \dots + |a_n| + |b_1| + \dots + |b_n|.$$

In (27) we can pass to an infinite sum on the right, obtaining

$$(28) \quad |c_1| + \cdots + |c_n| \leq \|a\| + \|b\|.$$

Since (28) holds for every n , we deduce that $\sum |c_i|$ converges and is bounded by the right side of (28). This proves that $c = \{c_i\}$ lies in l_1 and that $\|c\| \leq \|a\| + \|b\|$.

Example 7: Let m denote the set of all bounded sequences of real numbers, with the norm of $\{a_i\}$ given by $\sup |a_i|$. This is the infinite-dimensional analogue of 3(b) and 5(b).

Example 8: Let l_2 denote the set of all sequences $a = \{a_i\}$ of real numbers satisfying $\sum a_i^2 < \infty$. The norm $\|a\|$ is given by $\sqrt{\sum a_i^2}$. This example is called Hilbert space (more exactly, separable infinite-dimensional Hilbert space).

Again, the closure of l_2 under addition is not certain in advance. To settle this, and at the same time get the triangle inequality, both here and in Example 5(a), we need something called the *Schwartz inequality*. We derive it first for finite sums.

For $a = (a_1, \dots, a_n)$ and $b = (b_1, \dots, b_n)$ we introduce the inner product

$$(a, b) = a_1 b_1 + \cdots + a_n b_n.$$

Notice that (a, b) is symmetric, that it is linear in each factor, and that $(a, a) = \|a\|^2$, whence $(a, a) \geq 0$. Thus $(a + tb, a + tb) \geq 0$ for any real number t . Expanding this out, we get

$$(29) \quad (a, a) + 2t(a, b) + t^2(b, b) \geq 0.$$

Assume for the moment that $(b, b) \neq 0$. In (29) set $t = -(a, b)/(b, b)$. After simplifying we find

$$(a, b)^2 \leq (a, a)(b, b),$$

or, after extraction of square roots on both sides,

$$(30) \quad |(a, b)| \leq \|a\| \|b\|.$$

Now (30) obviously also holds when $b = 0$. Thus (30), the Schwartz inequality, is true unreservedly. We now find

$$\begin{aligned} \|a + b\|^2 &= (a + b, a + b) = \|a\|^2 + \|b\|^2 + 2(a, b) \\ &\leq \|a\|^2 + \|b\|^2 + 2\|a\| \|b\| = (\|a\| + \|b\|)^2. \end{aligned}$$

Hence $\|a + b\| \leq \|a\| + \|b\|$, the triangle inequality for Euclidean space, Example 5(a).

We can now make the transition to the infinite case. Let $a = \{a_i\}$, $b = \{b_i\}$ be given with $\sum a_i^2$ and $\sum b_i^2$ finite (the sum is over all positive integers i). The norms of a and b are the square roots of these sums. For any n , (30) is available and shows that

$$(31) \quad |a_1 b_1 + a_2 b_2 + \cdots + a_n b_n| \leq \|a\| \|b\|.$$

We let n go to infinity in (31) and see that the infinite sum $\sum a_i b_i$ converges and that moreover

$$|\sum a_i b_i| \leq \|a\| \|b\|.$$

We introduce the inner product $(a, b) = \sum a_i b_i$, and conclude the proof of the triangle inequality as above.

Example 9: The space of all continuous real functions on the closed interval $[0, 1]$ with

- (a) $\|f\| = \sup |f(x)|,$
- (b) $\|f\| = \int_0^1 |f(x)| dx,$
- (c) $\|f\| = \sqrt{\int_0^1 (f(x))^2 dx}.$

Example 10: Let Z be the integers and fix a prime p . Set $\|0\| = 0$ and, for $a \neq 0$, $\|a\| = 2^{-n}$ where n is the highest power of p dividing a . To verify the triangle inequality, we take any a and b in Z . If a or b or $a + b$ is 0, the checking of $\|a + b\| \leq \|a\| + \|b\|$ is trivial. Hence assume all three nonzero. We can write $a = p^m a_1$, $b = p^n b_1$ where m and n are non-negative integers and a_1, b_1 are prime to p . Then $\|a\| = 2^{-m}$ and $\|b\| = 2^{-n}$ by definition. Let us suppose, for definiteness, that $m \leq n$. Then $a + b$ is divisible by at least p^m . (*Remark:* If $m < n$, $a + b$ is exactly divisible by p^m , while if $m = n$, the power of p dividing $a + b$ may be higher.) Hence $\|a + b\| \leq 2^{-m}$. In short, $\|a + b\|$ is at most equal to $\max(\|a\|, \|b\|)$, a stronger statement than $\|a + b\| \leq \|a\| + \|b\|$. In fact, this metric space is ultrametric, as defined in Exercise 6 of Section 4.1.

Example 11: Let X be any set and let $C(X, R)$ be the set of all bounded real functions on X . For $f \in C(X, R)$, set $\|f\| = \sup_{x \in X} |f(x)|$. Example 7 is the special case of Example 11 where X is countably infinite, and Example 5(b) is the case where X is finite.

Example 12: We generalize Example 11 by having the functions go from X to an arbitrary metric space M instead of from X to the reals. For this purpose we have to reinterpret the word "bounded." The obvious way to do this is to call a function $f: X \rightarrow M$ bounded if its range has finite diameter. We write $C(X, M)$ for the space of bounded functions from X to M , taking as distance $D(f, g) = \sup D[f(x), g(x)]$. In Section 5.1 there is a discussion of the still more general case where X is also a metric space and the functions are restricted to be continuous.

Example 13: The following example has a geometric flavor. Generalized to higher-dimensional spaces, it is called *spherical geometry*. By identifying antipodally opposite points, one converts a spherical geometry to an *elliptic geometry*. This is one of the two non-Euclidean geometries;

the other (*hyperbolic geometry*) also gives rise to interesting metric spaces, but their definition is too complex to give here.

Let M be the circumference of the unit circle in the plane. In Cartesian coordinates, M is the set of all (x, y) with $x^2 + y^2 = 1$. Using complex numbers, we can describe M as the set of all $e^{i\theta}$, $0 \leq \theta < 2\pi$. As distance function between two points of M we take the shorter of the two arcs joining them. Analytically, if $0 \leq \theta \leq \phi < 2\pi$,

$$\begin{aligned} D(e^{i\theta}, e^{i\phi}) &= \phi - \theta && \text{if } \phi - \theta \leq \pi, \\ &= 2\pi - (\phi - \theta) && \text{if } \phi - \theta > \pi. \end{aligned}$$

2

Set Theory and Algebra

In this appendix we shall assemble, for the reader's convenience, brief definitions of some basic objects studied in modern algebra, and we shall then sketch several algebraic applications of transfinite methods.

A *group* is a set with a binary operation which is associative, has a two-sided unit element, and has the property that every element admits a two-sided inverse. If the operation is commutative, we call the group *abelian*. The group operation is normally written as multiplication, but if the group is abelian it is customary to use addition.

A *ring* is a set with two binary operations, addition and multiplication. Under addition a ring is an abelian group. Multiplication is assumed to be associative, and addition and multiplication are linked by the two distributive laws. A *field* is a ring in which the nonzero elements form a (multiplicative) abelian group. An *ideal* in a ring is a subset closed under addition and under multiplication by any element in the ring. (We shall discuss ideals only in commutative rings, and so we pass over the distinction between left, right, and two-sided ideals.)

A *vector space* V over a field K has the following elements of structure: V is an additive abelian group, and there is a multiplication between K and V , taking values in V . The axioms may be described as the list of natural laws that can be formulated. Elements u_i of V are *linearly inde-*

pendent if $\sum \lambda_i u_i = 0$ ($\lambda_i \in K$) implies that each $\lambda_i = 0$. A *basis* of V is a linearly independent subset of V which spans V (in the sense that every element of V is a—necessarily unique—linear combination of basis elements).

So much for definitions. We now sketch some applications of set theory to algebra.

(1) *Existence of a basis in a vector space V .* This is an application of Zorn's lemma which is a model of simplicity: Get a maximal linearly independent set in V by applying Zorn's lemma to the collection of all linearly independent subsets of V , partially ordered by inclusion.

(2) *Invariance of the number of elements in a basis.* Let the vector space V have two bases $\{u_i\}$, $\{v_j\}$, where i and j run over the index sets I and J . We are to prove that I and J have the same cardinal number. The case where either I or J is finite is a matter of basic linear algebra which we do not discuss (but we remark that it is very easy to rule out the possibility that one is finite and the other infinite). We shall assume that both I and J are infinite, and show that they have the same cardinal number. A typical element u_i has a unique expression in terms of the v 's, say

$$(32) \quad u_i = \alpha_1 v_{j_1} + \alpha_2 v_{j_2} + \cdots + \alpha_r v_{j_r}.$$

The order of the terms on the right of (32) is arbitrary but to be thought of as fixed. Let N be the set of positive integers. Let w be an extra symbol. We name a map f from $I \times N$ to $\{J, w\}$. On the portion of $I \times N$ with first coordinate i , where i is the index occurring in (32), f is defined by sending $(i, 1)$ to v_{j_1} , $(i, 2)$ to v_{j_2} , . . . , (i, r) to v_{j_r} , and (i, m) to w for $m > r$. We claim that f maps $I \times N$ onto $\{J, w\}$. Of course, w lies in the range. If some element of J is missing from the range of f , suppose that it corresponds to v_k . Write v_k in terms of the u 's, then these u 's in terms of the v 's. The result will be to write v_k as a linear combination of the remaining v 's, a contradiction of linear independence. Since $o(I \times N) = o(I)$ by Theorem 16, and the adjunction of w to J does not change its cardinal number, we deduce $o(I) \geq o(J)$. Of course the reverse inequality also holds, and so $o(I) = o(J)$.

Remark: The full force of Theorem 16 was not needed. The case where one of the cardinal numbers is \aleph_0 suffices and this is easy to deduce, for example, from Theorem 12.

(3) *Maximal ideals.* Let R be a commutative ring with unit. By a *maximal ideal* in R we mean an ideal $I \neq R$ such that there is no ideal properly between I and R . It is easy to see that any ideal (other than R) can be enlarged to a maximal ideal. While this is a routine application of Zorn's lemma, there is one point to be carefully noted. When we take the set-theoretic union of a chain of ideals, none of them equal to R , we have

to make sure that the union is not R . The presence of a unit element 1 does the trick, for none of the ideals of the chain contains 1, and hence the union does not contain 1.

(4) *Maximal subgroups.* In a ring without unit element, the argument just given does not apply, and in fact there need not exist maximal ideals. Likewise, a group need not have any maximal subgroups. But there is a variation on the theme that is worth recording: In a finitely generated group any proper subgroup can be enlarged to a maximal subgroup. Let G be generated by g_1, g_2, \dots, g_n and let H be a subgroup, $H \neq G$. For $i = 1, \dots, n - 1$ let H_i be the subgroup generated by H and g_1, \dots, g_i , and set $H_0 = H$. Take the largest j such that $H_j \neq G$. Then H_j and g_{j+1} generate G . As in our discussion of rings, we can enlarge H_j to a subgroup K maximal with respect to the exclusion of g_{j+1} . Evidently K is a maximal subgroup of G .

(5) *Free abelian groups.* In the language of direct sums—which are allowed to be infinite—a free abelian group is a direct sum of groups each of which is infinite cyclic (i.e., isomorphic to the additive group of all integers). For our purposes it will be convenient to describe a free abelian group G as one having a *basis* $\{u_i\}$ in the sense that every element of G has a unique expression as a linear combination of the u 's with integral coefficients. Here i runs over some index set I .

We shall now discuss the theorem that any subgroup of a free abelian group is free. On page 50 of Lefschetz's *Algebraic Topology* (*American Math. Society Colloquium Publ.* no. 27, 1942) it is asserted that for this theorem well-ordering gives a shorter, more intuitive proof than Zorn's lemma. I agree, although on page 44 of my *Infinite Abelian Groups* (Rev. ed., Univ. of Mich. Press, 1969) I have stubbornly given a Zorn style proof.

Here is a sketch of the proof in the style that uses well-ordering. As above, let G be a free abelian group, with the basis $\{u_i\}$ where i runs over the index set I . We now assume that I is well-ordered. Let g be any non-zero element in G . Then g has a unique expression of the form

$$(33) \quad g = a_1 u_{\alpha_1} + a_2 u_{\alpha_2} + \cdots + a_r u_{\alpha_r}$$

where a_1, \dots, a_r are integers and $\alpha_1, \dots, \alpha_r$ are members of I . It is a vital point that in (33) only a finite number of terms occur. There is thus a largest element among $\{\alpha_1, \dots, \alpha_r\}$, "largest" being meant in the ordering of I . Call this largest element the *index* of g .

We assume that a subgroup H of G is given, and we proceed to the proof that H is free. For any (momentarily fixed) β in I , we define a certain subset Z_β of the additive group Z of integers: It consists of 0, together with all coefficients of u_β that occur in elements of index β in H . We claim that Z_β is a subgroup of Z . If $m, n \in Z_\beta$ we have

$$h_1 = mu_\beta + \dots, \quad h_2 = nu_\beta + \dots$$

for suitable elements h_1, h_2 in H . We now verify that $m - n \in Z_\beta$. If $m = n$, there is no problem, for by definition 0 lies in Z_β . For $m \neq n$ the equation

$$(34) \quad h_1 - h_2 = (m - n)u_\beta + \dots$$

proves that $m - n \in Z_\beta$, for the remaining terms in (34) involve elements u_γ with $\gamma < \beta$. We now cite the elementary fact that any subgroup of Z is cyclic. Thus Z_β is cyclic. If $Z_\beta = 0$, there is no further activity concerning the index β . If $Z_\beta \neq 0$, we pick a generator r_β (it is natural to take r_β positive, but this is not essential). H contains an element of index β , in which u_β carries the coefficient r_β ; let such an element be v_β . (*Remark:* We have to pick v_β simultaneously for an infinite number of β 's; this is an amusing instance of how the axiom of choice can slip into an argument, almost unnoticed.) We claim that H is free, with the elements v_β constituting a basis.

To see this, we make an indirect argument, writing K for the subgroup generated by the v_β 's, and assuming that $K \neq H$. Among all the elements of H that are not in K we can pick one of minimal index (since I is well-ordered). Suppose that h' is that element and that γ is its index. Write

$$h' = pu_\gamma + \dots$$

Now p is an element of Z_γ and therefore is a multiple of r_γ , the generator of Z_γ . Let $h'' = h' - (p/r_\gamma)r_\gamma$. Then, since $v_\gamma \in K$, h'' is also an element of H lying outside K . Since the index of h'' is smaller than γ , we have our contradiction.

(6) *The Ulm invariants.* In advanced abelian group theory, there is an application of set theory which is remarkable in that both cardinal numbers and ordinal numbers play an essential role.

Fix a prime p . An abelian group G is called *primary* (for the prime p) if every element of G has order a power of p . Thus for any $x \in G$, there exists a positive integer n such that $p^n x = 0$. The power n is allowed to depend on x , and indeed the study of G is interesting only in the case where no fixed n works for all of G . If we form the series of subgroups

$$G \supset pG \supset p^2G \supset \dots \supset p^kG \supset \dots,$$

what we have just said may be rephrased as follows: In the interesting case this descending chain of subgroups does not reach 0 in a finite number of steps.

Cantor's basic idea can now be invoked: Why stop the chain at this point? We continue by defining $G_w = \bigcap p^k G$, $G_{w+1} = pG_w$, etc. We can go on to any ordinal number. In detail: If α is not a limit ordinal, we set $G_\alpha = pG_{\alpha-1}$; if α is a limit ordinal, we set $G_\alpha = \bigcap G_\beta$, the intersection

being taken over all $\beta < \alpha$. We cannot keep losing elements forever. Hence there must exist an ordinal λ at which stability sets in: $G_\lambda = pG_\lambda$. It turns out to be a harmless normalization to assume that $G_\lambda = 0$, and we do so.

Now we attach a cardinal number to each ordinal α less than λ . The way this is done needs motivation, which can be found in the various expositions of abelian groups which are available; we shall merely state it. Define P to be the subgroup of G consisting of all x with $px = 0$. The quotient group $(P \cap G_\alpha)/(P \cap G_{\alpha+1})$ can be regarded as a vector space over the field of integers mod p , and as such has a dimension. Call this cardinal number the α th *Ulm invariant* of G . Observe that the Ulm invariants constitute a function from ordinal numbers to cardinal numbers.

Ulm's theorem asserts that when G is countable, its Ulm invariants are a complete set of invariants. This decisive structure theorem is definitely false without the countability assumption. Recent years have seen a lot of progress in the study of uncountable primary abelian groups, but many mysteries remain.

3

The Transition to Topological Spaces

In this appendix topological spaces will be defined, and their relation to metric spaces will be briefly explored. It is hoped that this sketch may be helpful to a reader who plans to go on from metric spaces to general topology.

A *topological space* is a set X in which certain subsets have been singled out as distinguished and are called *open*. The following axioms are assumed:

- (a) The null set and the entire space X are open.
- (b) The union of any collection of open sets is open.
- (c) The intersection of a finite number of open sets is open.

A metric space becomes a topological space when we equip it with the sets defined to be open relative to the metric. There are, to be sure, many other examples of topological spaces, for otherwise the generalization from metric spaces to topological spaces would not be very exciting. We give an extreme example: It is possible that the only open sets are the null set and X . If X has at least two points, this topology cannot come from a metric.

We shall now survey, from the new point of view, the principal topics which were covered in our study of metric spaces.

(1) *Closed sets.* There is no problem in defining closed sets in a topological space: we simply take them to be the complements of open sets. Closed sets satisfy the dual properties: the null set and the whole space are closed, the intersection of any collection of closed sets is closed, the union of a finite number of closed sets is closed. If we wish, we can take closed sets as the primitive concept, and define open sets to be their complements.

(2) *Neighborhoods.* With the concept of “open” safely in hand, there is no problem about “neighborhood”. A neighborhood of x is a set containing an open set containing x .

(3) *Convergence.* The definition of convergence in the style that uses neighborhoods can be repeated verbatim: A sequence converges to a point if every neighborhood of the point contains the sequence, after a finite initial segment is deleted. But the concept loses nearly all of its significance. While it is true that any closed set contains the limits of its convergent sequences, the converse is false. Various generalizations of convergent sequences have been invented to cope with the difficulty: directed sets, Moore-Smith limits, filters. For nearly all purposes the difficulty can be avoided simply by not using convergence at all.

(4) *Continuity.* Let X and Y be topological spaces, f a function from X to Y . Continuity of f at a point of X can be phrased in terms of neighborhoods as in Theorem 38. Continuity of f in all of X , as in Theorem 40, asserts that f^{-1} carries open sets of Y into open sets of X .

(5) *Uniform continuity, uniform convergence, Cauchy sequence, completeness, precompactness.* None of these concepts make sense for topological spaces. There is, however, a generalization of metric spaces, due to A. Weil, which seems just right for the purpose. The spaces in question are called *uniform spaces*, or *spaces with a uniform structure*.

Another possible approach would be to stick to metrizable spaces and assert, for instance, that a given topological space is complete in all metrics, or complete in at least one metric. There exist some theorems formulated in this way, but they play a peripheral role in topology.

(6) *Homeomorphism versus isometry.* A *homeomorphism* between topological spaces X and Y is a mapping which is one-to-one, onto, and continuous both ways (compare Exercise 9 in Section 4.4). If X and Y are metric spaces, we have the companion notion of *isometry* to consider. An isometry is a homeomorphism, but the reverse is of course not true. An intermediate concept is to require uniform continuity both ways.

We call a property topological if it is invariant under homeomorphism. For example, compactness is topological and completeness is not.

(7) *Separability*. The three properties occurring in Theorem 59 become quite distinct for topological spaces, and so they get different names. In widely used (but not universally accepted) terminology, the existence of a countable dense subset maintains the designation *separable*; the reducibility of open coverings to countable coverings makes the space a *Lindelöf space*; and the existence of a countable open base is called the *second axiom of countability*. (The weaker “first axiom of countability” asserts that every point has a countable base of neighborhoods, in a certain natural sense.)

(8) *Compactness*. For ready reference, we assemble the three versions of compactness which were presented in Section 5.3 (Theorems 64 and 72).

- (a) Every sequence has a convergent subsequence.
- (b) Every infinite set has a limit point.
- (c) Every open covering admits a finite subcovering.

It is a fascinating chapter in the history of mathematics to trace how mathematicians slowly became aware that these three conditions are equivalent in metric spaces, but become quite distinct in general topological spaces.

Originally, (c) got the designation *bicompact*. In due course it was convincingly argued that (c) is the really important notion, and so it should get the brief name *compact*. Customary terminology today is to call (a) *sequential compactness* and (b) *countable compactness*.

(9) *Local compactness*. A topological space is *locally compact* if every point possesses a compact neighborhood. This important concept was not introduced in our account of metric spaces. We should not confuse local compactness with what we might call “bounded compactness” (every closed subset of finite diameter is compact). Bounded compactness is not a topological property. It implies local compactness but the reverse is not true.

(10) *Connectedness*. A topological space X is *connected* if the only subsets of X which are both open and closed are the null set and X . For example, it can be shown that the real line is connected, as is any interval on the real line. At the other extreme, a discrete space with more than one point is not connected. (Just as for metric spaces, a topological space is called *discrete* if every point—and hence every subset—is open.)

Connectedness was not mentioned in our account of metric spaces for the following reason: In the transition from metric spaces to topological spaces, not the slightest change occurs in the discussion of connectedness. Here the efficiency of covering the topic only once seems quite convincing.

(11) *Separation axioms*. A separation axiom in a topological space

asserts the existence of enough open sets of certain kinds. There is a big hierarchy of such axioms, and they have been intensively studied. We shall briefly note two separation axioms.

A *Hausdorff space* is a topological space with the following property: For any two distinct points x, y there exist disjoint open sets U, V with $x \in U, y \in V$. Here is a modest instance of how the Hausdorff axiom improves the behavior of topological spaces. Let X be a topological space, A a subset of X . There is a natural way in which A inherits a topology from X (one takes as the definition the property that occurs in Exercise 9 of Section 4.3). Suppose that A is compact. Easy examples show that this does not imply that A is closed. But if X is Hausdorff, then one can prove that the compact subsets of X are closed. That this is true in metric spaces was noted as a remark after Theorem 67; there we took the devious route of proving that compact subsets are complete and complete subsets are closed. This is a good example where revisiting and generalizing a topic, done for metric spaces, calls for a new nonsequential proof.

Sketch of the proof: We have A compact in the Hausdorff space X . Fix y in the complement of A . For $a \in A$ we have disjoint open sets U_a, V_a with $a \in U_a, y \in V_a$. A finite number of the U 's cover A . The intersection of the corresponding V 's is an open set which contains y and is disjoint from A .

Remark: Bourbaki advocates that "compact" should mean "compact Hausdorff," and he uses "quasi-compact" for compactness without Hausdorff. He seems to be winning this skirmish in the battle of nomenclature.

The second separation axiom that we shall mention is *normality*; it asserts that disjoint closed sets can be inserted into disjoint open sets. Its chief importance lies in the fact that one can deduce from normality the Tietze extension theorem, which asserts that a continuous real function on a closed subset of a normal space can be extended to be continuous on the whole space. As a consequence this is, in particular, true for metric spaces since metric spaces are normal (Exercise 20 in Section 4.3). If the Tietze theorem admitted an easier proof in the metric case, it would have been worth inserting in our account. But since the metric property does not seem to help, the Tietze theorem should be done in its natural setting of maximum generality.

One point of view on the separation axioms is that they head toward making the space resemble a metric space. In this connection, we mention that a decisive characterization of metric spaces was given independently by Nagata and Smirnov. It is somewhat too complex to state here, and we invite the reader to consult one of the available treatises on general topology.

(12) *Products.* The subject of Cartesian products is a weak spot of metric space theory and a triumph for the generalization to topological spaces.

As we noted in Section 6.1, already for a finite product of metric spaces we are in doubt as to the best choice of a metric. The situation for an infinite product is even less attractive.

For a product of (any number of) topological spaces, a unique natural topology awaits us: It is defined by taking as a base for open sets all products of open sets in the factors, subject to the restriction that all but a finite number of the constituent open sets are required to be the whole space.

A key result is Tychonoff's theorem: *The Cartesian product of compact spaces is compact.*

Metric considerations fit in best as the following theorem: *The Cartesian product of topological spaces, each having at least two points, is metrizable if and only if each factor is metrizable and the number of factors is countable.*

(13) *Fixed points and category.* The material of Section 6.2 (fixed points) has stubbornly resisted a convincing generalization to topological spaces. Until recently the same was true for Section 6.3 (category), but there have been new developments. The paper *Cocompactness* by Aarts, de Groot, and McDowell in the 1970 *Nieuw Arch. Wisk.* is a pertinent reference.

Selected Bibliography

- Abbott, J. C. *Sets, Lattices, and Boolean Algebras*. Boston: Allyn and Bacon, Inc., 1969. Includes cardinal and ordinal numbers, transfinite methods, and a discussion of axiomatic set theory.
- Abian, A. *The Theory of Sets and Transfinite Arithmetic*. Philadelphia: Saunders, 1965. In addition to the standard transfinite topics, includes a construction of the real numbers.
- Alexandroff, P. S. *Einführung in die Mengenlehre und die Theorie der Reellen Funktionen*. 4th ed. Berlin: 1967. (Translated from Russian, originally published in 1948.) After a brief treatment of set theory, the book proceeds to point sets in Euclidean space and then to metric spaces. Appendixes treat topological spaces.
- Bachmann, H. *Transfinite Zahlen*. 2nd ed. New York: Ergebnisse der Math., Springer, 1967. Along with Sierpinski, this is a standard comprehensive treatise on the subject.
- Bourbaki, N. *Theory of Sets*. Reading, Mass.: Addison-Wesley, 1968. In one volume and in English, this contains Bourbaki's treatment of set theory. Like all of Bourbaki's works, it is a major contribution to the literature. The historical note (pp. 296-346) is especially worthy of attention.

- _____. *General Topology*. 2 vols. Reading, Mass.: Addison-Wesley, 1966. With Bourbaki's famous logical plan clicking to perfection, general topological spaces and uniform spaces precede metric spaces by several light years.
- Brown, R. *Elements of Modern Topology*. New York: McGraw-Hill, 1968. The first half of the book moves briskly through the standard topics, with topological spaces getting priority. The second half is an introduction to homotopy theory, featuring the fundamental groupoid.
- Bushaw, D. *Elements of General Topology*. New York: John Wiley and Sons, 1963. An "anti-metric" version. Uniform spaces are included.
- Copson, E. T. *Metric Spaces*. Cambridge: Cambridge University Press, 1968. Metric spaces only. Chapter 8, on applications to analysis, is unusually extensive.
- Dieudonné, J. *Foundations of Modern Analysis*. 2 vols. New York: Academic Press, enlarged and corrected printing 1969–70. An excellent treatment of metric spaces is given, leading up to his coordinate-free treatment of calculus.
- Dugundji, J. *Topology*. Boston: Allyn and Bacon, 1965. There is a fair-sized introduction to set theory in the first two chapters. Metric spaces make their entrance late (Chapter 9) but are then treated in detail.
- Fairchild, W. W., and C. Ionescu Tulcea. *Sets*. Philadelphia: Saunders, 1970. A short paperback in which the set theory is essentially naive, and the transfinite portion does not go very far. An unusual feature is a final chapter on combinatorial analysis.
- Fraenkel, A. A. *Abstract Set Theory*. North-Holland, 1953. A very readable discursive account, with occasional reference to axioms.
- _____. *Einleitung in die Mengenlehre*. 3rd ed. New York: Springer, 1928. Reprinted by Dover, 1946. A classical account, with unusually extensive historical and philosophical material.
- Gemignani, M. C. *Elementary Topology*. Reading, Mass.: Addison-Wesley, 1967. Some metric spaces first. There is a concluding chapter on homotopy.
- Halmos, P. R. *Naive Set Theory*. New York: Van Nostrand, 1960. The fundamental ideas are presented in a sprightly style, in the spirit of axiomatic set theory.
- Hausdorff, F. *Mengenlehre*. Leipzig. 1st ed. 1914, 2nd ed. 1927, 3rd ed. 1937. 1st ed. reprinted by Chelsea, 1949; 3rd ed. by Dover, 1944; 3rd ed. translated by J. Aumann et al. and published by Chelsea, 1957, under the title *Set Theory*. The full title of the 1st ed. is *Grundzüge der Mengenlehre*. The first edition is an imperishable classic: An extensive treatment of set theory is followed by the first general treatment of topological spaces (with the Hausdorff separation axiom built in). In the second edition he changed his mind, omitting a good deal of the set theory and giving metric spaces the place of honor. In the third edition a small amount of additional material was added.

- Hayden, S. and J. F. Kennison. *Zermelo-Fraenkel Set Theory*. Columbus, Ohio: Charles Merrill, 1968. The standard topics are discussed on the basis of axiomatic set theory in the style of Zermelo-Fraenkel. An appendix discusses other axiom systems briefly. A construction of the real numbers is included.
- Hocking, J. G. and G. S. Young. *Topology*. Reading, Mass.: Addison-Wesley, 1961. Metric spaces have a place of honor but not the first place. A large amount of algebraic topology is covered.
- Kamke, E. *Theory of Sets*. Trans. by F. Bagemihl. New York: Dover, 1950. A brief account in naive style.
- Kelley, J. L. *General Topology*. New York: Van Nostrand, 1955. In this widely used treatise metric spaces have to wait until page 118.
- Mendelson, B. *Introduction to Topology*. Boston: Allyn and Bacon, 1962. An extensive chapter on metric spaces precedes the general topological space.
- Monk, J. D. *Introduction to Set Theory*. New York: McGraw-Hill, 1969. The basis is axiomatic; the treatment of ordinal and cardinal numbers goes fairly far.
- Moore, T. O. *Elementary General Topology*. Englewood Cliffs, N.J.: Prentice-Hall, 1964. The general topological space is treated at once, but there is ample attention to metric spaces.
- Newman, M. H. A. *Elements of the Topology of Plane Sets of Points*. 2nd ed. Cambridge: Cambridge University Press, 1961. As the title indicates, special results in the plane are the main objective. However, the book furnishes as well an excellent introduction to metric spaces.
- Rotman, R. and G. T. Kneebone. *The Theory of Sets and Transfinite Numbers*. New York: Elsevier, 1966. The Zermelo-Fraenkel axioms are presented, but the text is accessible without them. In a relatively brief space, the standard topics are well covered.
- Sierpinski, W. *Cardinal and Ordinal Numbers*. New York: Hafner Publishing Co., 1958. Along with Bachmann, this is an authoritative comprehensive treatise on infinite numbers.
- Sigler, L. E. *Exercises in Set Theory*. New York: Van Nostrand, 1966. A collection of exercises, designed principally to accompany Halmos's *Naive Set Theory*.
- Simmons, G. F. *Introduction to Topology and Modern Analysis*. New York: McGraw-Hill, 1963. A full chapter on metric spaces precedes the general topological space.

Index

- Aarts, J., 131
- Abelian group, 122
 - free, 124
 - normed, 117
 - primary, 125
- Algebraic number, 23
- Antisymmetry, 19, 32
- Axiom of choice, 21, 36, 58
 - countable, 21, 64, 76, 80
- Axioms of countability, 129

- Belluce, L., 111
- Benacerraf, P., 65
- Bernstein, F., 33
- Birkhoff, G., 33
- Boehm, G., 65
- Boole, G., 13
- Boolean algebra, 13
- Borel, F., 59
- Bottom element, 11
- Bound
 - greatest lower, 11, 67
 - least upper, 11, 67
 - lower, 11
 - upper, 11, 37
- Burali-Forte, C., 56
 - paradox, 56

- Cantor, G., 24, 28, 29, 30, 33, 36, 49, 57, 64, 125
- Cardinal
 - addition, 40
 - exponentiation, 44
 - limit, 57
 - multiplication, 41
 - number, 27
- Cartesian product, 19, 41, 106, 131
- Category, 111, 131
- Cauchy sequence, 84
- Chain, 10
- Class, 1
- Closed
 - ball, 70, 73
 - mapping, 83

- Closed (*Cont.*)
 - set, 75, 128
- Closure, 77
- Cofinal, 26, 55
- Cofinite, 9
- Cohen, P. J., 65
- Collection, 1
- Compactness, 99, 129
 - bounded, 129
 - countable, 129
 - local, 129
 - sequential, 129
- Complement, 6
- Complemented lattice, 12
- Complete
 - conditionally, 14, 27, 67
 - lattice, 12
 - metric space, 86, 99, 128
- Completion, 90, 101
- Composite function, 15
- Condensation point, 96
- Connected, 75, 129
- Continuous, 80, 128
 - uniformly, 88, 101, 128
- Continuum hypothesis, 47, 64
 - generalized, 66
- Convergence, 75, 128
 - uniform, 90
- Convex
 - metric space, 94
 - subset of a chain, 14
- Countable, 22
- Countably infinite, 22
- Cycle, 33
 - bilateral, 34
 - unilateral, 34

- Dedekind, R., 33
- De Groot, J., 131
- Dense, 81
 - nowhere, 111
- Descending chain condition, 50, 58
- Diameter, 69, 99
 - finite, 69, 99
- Discrete, 74
- Domain, 14
- Duality, 7

- Empty set, 2
- Equality (of sets), 2
- Euclid, 69
- Euclidean
 - plane, 68
 - space, 86, 100

- Field, 122
- Function, 14
 - characteristic, 18, 44
 - identity, 15
 - product (or composite), 15

- Gelfond, A., 25
- Geometry
 - elliptic, 120
 - hyperbolic, 121
 - spherical, 120
- Gödel, K., 65
- Grelling, K., 30

- Halmos, P., 61
- Hausdorff, F., 36, 60
- Hausdorff space, 130
- Hermite, C., 25
- Hilbert, D., 25, 59, 64
- Hilbert space, 93, 119
- Homeomorphism, 83, 128

- Ideal, 50, 122
 - maximal, 123
- Identity function, 15
- Image, 14
 - inverse, 16, 81
- Inclusion, 2
- Index set, 8
- Intersection, 2, 4
- Isolated, 74
- Isometry, 69, 104, 128
- Isomorphism, 16
 - order, 16

- Jourdain, P., 28

- Kirk, W., 111
- König, J., 66

- Lang, S., 110
 Lattice, 12
 complete, 12
 conditionally complete, 14, 27, 67
 distributive, 12
 Lebesgue
 integral, 92
 number, 103, 105
 Lefschetz, S., 124
 Limit, 75
 point, 77, 98
 Lindelöf space, 129
 Lindemann, F., 25
 Lipschitz condition, 109

 Mac Lane, S., 33
 Map (or mapping), 14
 Maximal element, 37
 Mazur, S., 112
 McDowell, R., 131
 Metric space, 67

 Nagata, J., 130
 Neighborhood, 73, 128
 Normal, 130

 One-to-one, 14
 correspondence, 15, 28
 Onto, 14
 Open
 ball, 70
 base, 94
 covering, 95
 mapping, 83
 set, 70, 127
 Ordered pair, 19
 Ordinal number, 28, 55
 countable, 56
 initial, 56
 limit, 50

 Partially ordered set, 9
 Periodic, 17
 locally, 17
 Power set, 11, 29

 Precompact, 103
 Product
 function, 15
 metric, 106
 Putnam, H., 65

 Quasi-compact, 130
 Quine, W., 30

 Range, 14
 Reflexive, 19, 32
 Relation, 19
 equivalence, 19
 Ring, 122
 Robinson, R., 25
 Russell, B., 30

 Schmidt, E., 59
 Schneider, T., 25
 Schröder, E., 33
 Schwartz inequality, 119
 Segment, 50
 lower, 14
 Separable, 94, 100, 129
 Separation axioms, 129
 Set, 1
 empty (or null), 2
 index, 8
 partially ordered, 9
 power, 3, 12, 19
 Smirnov, Y., 130
 Smullyan, R., 65
 Sphere, 78
 Stone, M. H., 6
 Subcovering, 95, 101
 Subsequence, 85
 Subset, 3
 Subtraction (of sets), 7
 Symmetric, 19
 difference, 6, 9

 Tietze theorem, 130
 Top element, 11
 Topological space, 127
 Totally bounded, 103
 Totally unordered, 10

- Tower, 61
Transcendental number, 24
Transfinite
 construction, 57
 induction, 57
Transitive, 3, 19, 32
Triangle inequality, 68, 117
Tychonoff theorem, 131

Ulm invariants, 124
Ultra-metric, 70, 120
Uncountable, 21, 24

Uniform space, 128
Union, 4

Vector space, 122
Venn, J., 5
Venn diagram, 5
Von Neumann, J., 31

Weil, A., 128
Well-ordered, 36, 49
Whitehead, A. N., 30

Zermelo, E., 36, 58
Zorn, M., 36
Zorn's lemma, 9, 36, 52, 58, 123